

Jia-Jia (Jay) Lin
CSE 415, Winter 2020
February 3rd, 2020

Assignment 4: Worksheet

1. Blind Search

- a. A possible state representation is a 2D int array with $6 \times 6 = 36$ total elements representing each square in Sudoku.
- b. `def successors(s):`
 `list = []`
 `# loop through all possible choices from this state`
 `for i in range(7):`
 `for square in s state:`
 `if square is Empty: list.append(state with i in Z)`
 `return list`
- c. `def is_goal(s):`
 `return (s == 2D array of final correct solution)`

2. Heuristic Search

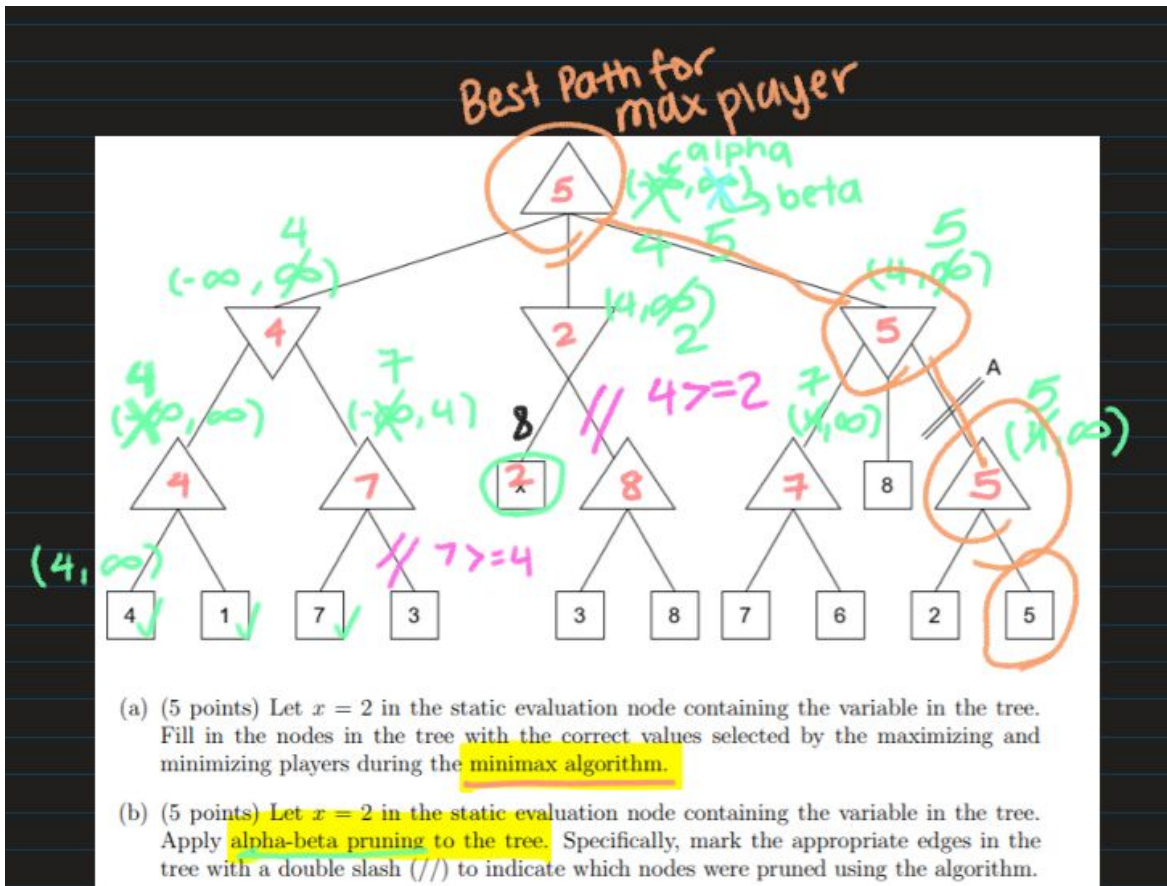
- a. When selecting a heuristic, the challenge is deciding what the value represents, making sure it is reasonable, admissible, and/or consistent. Trade-offs include making the heuristic values too small (low risk of overestimating the distance to the goal node but not very efficient compared to a regular heuristic search) or too big (very efficient but could overshoot the goal node).
- b. h_1 and h_2 are admissible.
 - i. h_3 has a heuristic value that could potentially overshoot the actual distance to the goal from Node B ($12 > 11$)
- c. None of the heuristics shown above are consistent
 - i. h_1 : s_0 to Node B, $14 - 4 = 10$, which is not ≤ 3 , the cost.
 h_2 : $14 - 10 = 4$ not ≤ 3 , h_3 : $12 - 7 = 5$ not ≤ 3 from Node B to D
 - ii. The above shows that the condition to be consistent, that every node's heuristic value being less than or equal than its successors' heuristic + cost to get there, is not met.
- d. An A* search algorithm is optimal when using a consistent and admissible heuristic. While none of them are consistent, I would use h_2 because it does not overestimate the distance to the goal node, but still has a better efficiency than h_1 .
- e. Final Solution Path: $[(s_0, 6), (B, 8), (D, 9), (H, 10), (E, 12), (\text{gamma}, 14)]$
 - i. Path A* traversed: $s_0 \rightarrow B \rightarrow A \rightarrow G \rightarrow D \rightarrow F \rightarrow H \rightarrow E$

Table for Question 2 Part E: A* Search path

Current Node	OPEN list	CLOSED list
s0	[(B,8),(A,9),(G,9),(F,10), (D,11)]	[s0, 6]
B	[(A,9),(G,9), (D,9), (F,10),(C,17)]	[(s0,6),(B,8)]
A	[(G,9), (D,9), (F,10),(C,17)]	[(s0,6),(B,8)]
G	[(D,9), (F,10),(H,15),(C,17)]	[(s0,6),(B,8), (G,9)]
D	[(F,10),(H,10),(C,17)]	[(s0,6),(B,8), (D,9)]
F	[(H,10),(C,17)]	[(s0,6),(B,8),(D,9)]
H	[(E,12),(C,17),(gamma, 20)]	[(s0,6),(B,8), (D,9), (H,10)]
E	[(gamma,20), (C,17)]	[(s0,6),(B,8), (D,9),(F,10),(H,10),(E,1 2)]
gamma	[(C,17)]	[(s0,6),(B,8), (D,9), (H,10),(E,12),(gamma, 14)]

3. Adversarial Search (note: red pen = minimax, green pen = alpha-beta pruning)

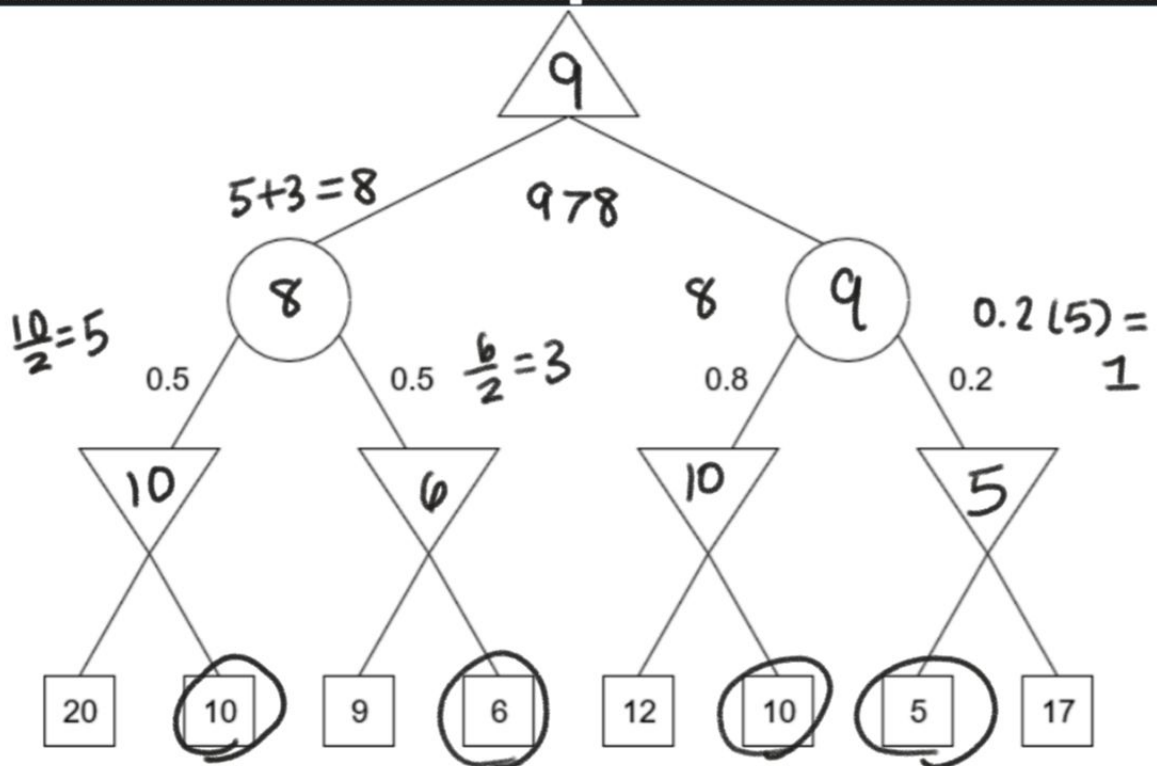
a. (part a and b drawing shown below)



c. What is the smallest value of x where A is pruned using alpha-beta?

$x = 7$, because the minimizing player will choose 7 in the middle branch, which then gets passed to the right side of the tree later. At the right child, it will be $[4, 7]$ and goes to its left child first, which is the maximizing player (chooses 7) $\rightarrow [7, 7]$. 7 is greater than or equal to 7, so the remaining branches are pruned.

d. Consider the following expectimax game tree. (Answers in drawing below)



4. Markov Decision Processes

a. Actions and states

- i. Actions for this MDP: $a_1 = \text{Roll } 1$, $a_2 = 3$, $a_3 = 5$, or $a_4 = 6$, or $a_0 = \text{STOP}$
- ii. 8 total states: $s_0 = \text{Initial state}$, sums = $\{s_1, s_2, s_3, s_4, s_5, s_6\}$, $s_7 = \text{Final State}$

b. Full transition function

- i. $T(s, a, s') = 0.25$, probability of rolling a number on a tetrahedral dice; given $s = s_0$
 1. a is within $\{a_1, a_2, a_3, \dots\}$ and NOT a_0
 2. s' is within $\{s_1, s_2, s_3, \dots\}$ and NOT s_7
- ii. Example: $T(s_0, a_1, s_1) = 0.25$

c. Reward functions

- i. $R(s_0, a_2, s_3) = 3$ (total sum of 3)
- ii. $R(s_0, a_3, s_5) = 5$
- iii. $R(s_0, a_4, s_6) = 6$
- iv. $R(s_1, a_1, s_2) = 2$
- v. $R(s_1, a_2, s_4) = 4$
- vi. $R(s_1, a_3, s_6) = 6$
- vii. **$R(s, a, s') = -10000$ where $s' = s_7$ and both $s = s_0$ and $a = a_0$, OR a is not a_0**
 1. In words: The reward is very bad if the player immediately forfeited the game at the first turn or if they got a bust.
 2. Example: $R(s_0, a_0, s_7) = -10000$
- viii. **$R(s, a, s') = 1$, where $s' = s_1$**
 1. Example: $R(s_0, a_1, s_1) = 1$
- ix. **$R(s, a, s') = 2$, where $s' = s_2$**
 1. Example: $R(s_1, a_1, s_2) = 2$
- x. $R(s, a, s') = 3$, where $s' = s_3$
- xi. $R(s, a, s') = 4$, where $s' = s_4$
- xii. $R(s, a, s') = 5$, where $s' = s_5$
- xiii. **$R(s, a, s') = 10000$, where $s' = s_6$**

1. The maximum / best possible score

- d. The optimal policy is one where the agent plays carefully but reasonably: avoids stopping the game too early and not overshooting the number 7. The agent should have an increasingly higher preference for stopping when reaching a total sum close to 7.

5. Computing MDP State Values and Q-Values

- 0
- 0
- 8.5
- 4.6
- 34.4
- 17.4

(see work for value iteration below)

s	a	s'	T(s,a,s')	R(s,a,s')
s ₁	x	s ₁	1	0
s ₁	x	s ₂	0	0
s ₁	y	s ₁	0.5	4
s ₁	y	s ₂	0.5	1
s ₁	z	s ₁	0.4	0
s ₁	z	s ₂	0.6	10
s ₂	x	s ₁	0	0
s ₂	x	s ₂	1	0
s ₂	y	s ₁	0.9	2
s ₂	y	s ₂	0.1	8
s ₂	z	s ₁	1	1
s ₂	z	s ₂	0	0

states s₁ and s₂.

V₀ = reward value at step 0

(s₁) at s₁

✓ x, s₁ = 0

✓ x, s₂ = 1

s₁ and s₂.

✓ y, s₁ = 0.9 [2 + 1.8] = 1.8

✓ y, s₂ = 0.1 [8 + 1.8] = 0.8

✓ z, s₁ = 1 [1] = 1

✓ z, s₂ = 0

V₂(s₂):

1 [0 + 4.6] = 4.6

0.9 [2 + 4.6] = 5.94

0.1 [8 + 4.6] = 1.26

1 [1 + 4.6] = 5.6

γ = 1
s₁ → x
V₀(s₁) = 1 · 0 = 0

s₁ → y
V₁(s₁) = max_a ∑_{s'} T(s,a,s') [R(s,a,s') + γ V₀(s')]

✓ x, s₁: 1 [0 + 0] = 0

✓ x, s₂: 0

✓ y, s₁: 0.5 [4 + 0] = 2

✓ y, s₂: 0.5 [1 + 0] = 1/2

✓ z, s₁: 0.4 [0 + 0] = 0

✓ z, s₂: 0.6 [10 + 0] = 6

V₂(s₁) = max_a ∑_{s'} T(s,a,s') [R(s,a,s') + γ V₁(s')]

1 [0 + 1.8] = 1.8

0.5 [4 + 1.8] = 2.9

0.5 [1 + 1.8] = 0.95

0.4 [0 + 1.8] = 0.72

0.6 [10 + 1.8] = 6.48

V₂(s₁) = 8.5

0.5 [4 + 8.5] = 6.25

0.5 [1 + 8.5] = 4.75

0.4 [0 + 8.5] = 3.4

0.6 [10 + 8.5] = 11.1

V₂(s₁) = 34.4

2 + 1/2 + 6 = 8.5

k=1

8.5

V₂(s₁) = 34.4

- g. (Question 5 continued) 8.5
- h. 11.0
- i. 14.5
- j. 4.6
- k. 7.2
- l. 5.6

(see work for Q values below)

s	a	s'	T(s, a, s')	R(s, a, s')
s ₁	x	s ₁	1	0
s ₂	x	s ₂	0	0
s ₁	y	s ₁	0.5	4
s ₁	y	s ₂	0.5	1
s ₁	z	s ₁	0.4	0
s ₁	z	s ₂	0.6	10
s ₂	x	s ₁	0	0
s ₂	x	s ₂	1	0
s ₂	y	s ₁	0.9	2
s ₂	y	s ₂	0.1	8
s ₂	z	s ₁	1	1
s ₂	z	s ₂	0	0

take x action, move another time:

$Q_2(s_1, x) = 0$
 $Q_2(s_1, x) = V_1(s_1) = 8.5$
 s_1 to anything

$Q_2(s_1, y)$
 $Q_1(s_1, y) = 0.5(4) + 0.5(1) = 2.5$
 $Q_2(s_1, y) = 2.5 + 8.5 = 11.0$

$Q_2(s_1, z)$:
 $Q_1(s_1, z) = 0 + 0.6(10) = 6$
 $Q_2(s_1, z) = 6 + 8.5 = 14.5$

$Q_2(s_2, x)$
 $Q_1(s_2, x) = 0$
 $Q_2(s_2, x) = 0 + 1 = 1$

$Q_2(s_2, y)$:
 $Q_1(s_2, y) = 0.9(2) + 0.1(8) = 2.6$
 $Q_2(s_2, y) = 2.6 + 4.6 = 7.2$

$Q_2(s_2, z)$:
 $Q_1(s_2, z) = 1$
 $Q_2(s_2, z) = 1 + 4.6 = 5.6$

Finally, compute V_2 and Q_2 for both states, but with $\gamma = 0.5$.

- m. 29.75
- n. 15.1
- o. 8.5
- p. 11.0
- q. 14.5
- r. 4.6
- s. 7.2
- t. 5.6

(see work below)

- Reasoning for Q-values
 - Since $Q_0 = 0$, the Q_2 values didn't change regardless of gamma value because it is $(\gamma \cdot 0 + \text{the reward}) \cdot \text{probability}$ for the first step and then add $V_1(s)$ which still involved multiplying gamma by 0 since $V_0 = 0$

Finally, compute V_2 and Q_2 for both states, but with $\gamma = 0.5$.

(m). $V_2(s_1) = \underline{29.75}$
 (n). $V_2(s_2) = \underline{15.1}$
 $V_2(s_2) | \gamma=1 = 4.6$

$\gamma=1.2$
 $V_1(s_1) = 8.5$
 because $V_0=0$

$4.25 + 4.125 + 2.625 + 1.7$
 $+ 8.55 = 21.25$

$2.3 + 3.87 + 1.03 + 3.3 = 10.5$

(o). $Q_2(s_1, x) = \underline{8.5}$
 (p). $Q_2(s_1, y) = \underline{11.0}$
 (q). $Q_2(s_1, z) = \underline{14.5}$
 (r). $Q_2(s_2, x) = \underline{4.6}$
 (s). $Q_2(s_2, y) = \underline{7.2}$
 (t). $Q_2(s_2, z) = \underline{5.6}$