

CSE150B Spring 2024 Final Exam

Please, finish reading everything on this page before you start.

- Open book, web, and everything. No chatgpt. Do all the work **completely by yourself** without any discussion with others. Any form of communication related to the course before the submission deadline will be considered academic integrity violation.
- For clarification questions, ask me via **direct messages** on Slack (@Sean) and do **not** post anything in public channels. If any general correction/clarification is needed, I'll make announcements in the #general channel on Slack.
- Be assured that I will read every message, and if you do not see my reply within an hour or so (when my Slack status is green), it means my reply is "no comment" because the answer will be clear after re-reading the questions or going through class materials.
- Submit the following:
 1. Submit photos of your work (assignment: "Final").
 2. Submit a video of your work at the following link:
<https://www.dropbox.com/submit/video?ref=link>
- **Make sure to** let the camera know you are working. Your work is spotty. You need to show your face and turn off camera and only show the answer screen.
- There is no strict time limit. Do not make it too long (20 minutes) as for each question you should be more thoughtful. Think about the answer before you answer. Think about what the grader/aminer about what you understand, while you are answering.
- The grading will be based on your answers. To ensure academic integrity, the grader will not only ask you to explain your answer but also help get to the bottom of it. If you ask something, you need to explain it. Otherwise that question will be considered unanswered.
- Make sure the file you upload is a video. I don't know if Dropbox sends you confirmation emails. Look for a link at the bottom of the page indicating that you have successfully uploaded your files. You can upload multiple times, but we do not guarantee that we will only check the last one. So try your best to send only a final version of the video.
- Zoom is highly recommended for doing your recording. If you recorded using your phone or in some other ways, try compressing the video into mp4 format (can use relatively low resolution) and avoid making the file too large. If you find yourself uploading gigabytes of video, you should change the file.
- Your overall letter grade for the course will only depend on the sum of all raw points you obtained in the class, and there is no notion of "percentage" in this class. You only need to get enough points for the threshold of the grade you want.

1 Classical Search

Question 1 (7 Points). The undirected graph in Figure 1-(Left) has four nodes (A,B,C,D). The number on each edge indicates the cost of going through the edge – e.g., going from A to B has a cost of 5, same as going from B to A. Set node A as the start, and D the goal node. Answer the following questions and explain your reasoning.

1. In the correct algorithm for Dijkstra, if we visit a node that is already in the frontier, we need to compare the cost of this node from the current path with its cost in the frontier (the last two lines in the pseudocode of Slide 36 of this chapter 1 Classical Search.pdf). Now, on this given graph, if we do not do that comparison and replacement (i.e., remove the last line in the pseudocode), then what is the path from A to D that this wrong modification of Dijkstra algorithm will return? Explain the steps of how you obtained it.
2. What is the path that will be returned by the correct Dijkstra algorithm instead?
3. In your PA1, we agreed that it was ok to not do this frontier check and replacement, and still always get the optimal path. Explain why it was ok in the graphs in PA1. In your explanation, you can use the fact that the Dijkstra algorithm always returns the lowest path cost.
4. Define the following heuristic h : $h(A) = 1, h(C) = 7, h(D) = 0$. Is h a consistent heuristic?
5. What is the path from A to D that this heuristic h will return?

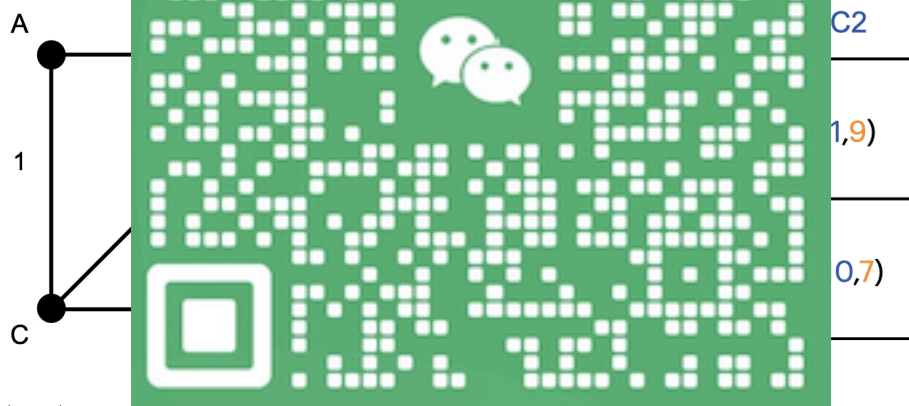


Figure 1: (Left) Graph for the classical search question. (Right) Game for adversarial search question.

2 Adversarial Search

Question 2 (3 Points). Consider the game in Figure 1-(Right). There are two players: the column player and the row player. The column player goes first, and chooses either the first column (C1) or the second (C2), and after that, the row player chooses either the first row (R1) or the second (R2). After these two steps, the game ends and the pair of two numbers located by these two choices is the outcome of the game. In each pair of numbers, the first number is the reward that the column player will receive, and the second number is what the row player will receive. For instance, if the column player chooses C1, and then the row player chooses R1, then the outcome is (10,2), which means the column player will receive a reward of 10 and the row player a reward of 2.

Assume that both players are rational players and want to maximize their own reward (no minimizers), what should each of their action be, and what is the final outcome of the game under those actions?

In your answer, draw the game tree that is similar to the minimax tree, just that now the two players are no longer the max and min players (just annotate their nodes as column and row players). Explain your reasoning on the tree. I hope you realize this is how the algorithms we talked about in minimax can be applied to games where players are not directly adversarial to each other (not zero-sum).

3 Markov Decision Processes and Reinforcement Learning

Consider the following MDP:

- State space: $S = \{s_1, s_2\}$
- Actions: $A(s_1) = \{a_1, a_2\}$ and $A(s_2) = \{a_1, a_2\}$.
- Transition model:

- $P(s_1|s_1, a_1) = 0.2$
- $P(s_1|s_2, a_1) = 0.7$

- Rewards: $R(s_1) = 0$

- Discount factor: $\gamma = 0.9$

Question 3 (2 Points). Draw the MDP graph. Show the states, actions, and the appropriate edges in the graph.

Question 4 (3 Points). Suppose $V_0(s_1) = V_0(s_2) = 0$. That is, perform two steps of value iteration (i.e., $V_{k+1}(s) = \sum_a P(s'|s, a) [R(s') + \gamma V_k(s')]$). Note that each V_i is a vector of the values for the two states. How would you compute V_1 and V_2 , where each V_i is a vector of the values for the two states?

Question 5 (2 Points). Suppose you know what the optimal values $V^* = (V^*(s_1), V^*(s_2))$ are, within 0.1 error from the optimal values. That is, $|V_k(s_i) - V^*(s_i)| \leq 0.1$ for all i . and we want $|V_k(s_1) - V^*(s_1)| \leq 0.01$ and $|V_k(s_2) - V^*(s_2)| \leq 0.01$. How many iterations does k need to be to converge to under this tolerance? What is the computation involved.

Question 6 (2 Points). Suppose you know what the optimal values $V^*_{M_1}(s_1), V^*_{M_1}(s_2)$ to represent the optimal values for M_1 are. Write $V^*_{M_2}(s_1), V^*_{M_2}(s_2)$ to represent the optimal values for M_2 in terms of $V^*_{M_1}(s_1), V^*_{M_1}(s_2)$. that the rewards are now $r_{M_2}(s_1) = 2$ and $r_{M_2}(s_2) = -4$. Namely, M_2 differs from M_1 only in that the reward on each state is twice as large. We write the optimal values of M_2 as $V^*_{M_2}(s_1)$ and $V^*_{M_2}(s_2)$.

Now, suppose you know what $V^*_{M_1}(s_1)$ and $V^*_{M_1}(s_2)$ are (which you don't, and you do not need to implement it to find out), then what should $V^*_{M_2}(s_1)$ and $V^*_{M_2}(s_2)$ be? That is, write down the values of $V^*_{M_2}(s_1)$ and $V^*_{M_2}(s_2)$ using the values of $V^*_{M_1}(s_1)$ and $V^*_{M_1}(s_2)$. Explain your reasoning.

Question 7 (4 Points). We now perform Q -learning in the MDP defined in the beginning, i.e., M_1 , without accessing the transition probabilities. We initialize all Q -values to be zero $Q(s_i, a_i) = 0$. For simplicity we will use a fixed learning rate parameter $\alpha = 0.1$ in all temporal-difference updates (so it does not change with the number of visits). We start from the initial state s_2 , and suppose we perform/observe the following three steps in the learning process in the MDP:

- (Step 1) At the initial state s_2 , we take the action a_1 , and then observe that we are transitioned back to s_2 by the MDP.

- (Step 2) We then take the action a_1 again from s_2 , and observe that we are transitioned now to s_1 by the MDP.

- (Step 3) Now at the state s_1 , we take the action of a_1 yet again, and then observe that we are transitioned to state s_2 by the MDP.

After each step of these three steps, we perform temporal-difference updates for the Q -values on the relevant state-action pairs. Now answer the following questions.

1. What values do we have for $Q(s_1, a_1)$ and $Q(s_2, a_1)$ now, after these three steps of updates? Write down how you obtained them.
2. Suppose from here we will use the ε -greedy strategy with $\varepsilon = 0.3$, which means that with ε probability we will use an arbitrary action (each of the two actions will be chosen equally likely in this case), and with $1 - \varepsilon$ probability we will choose the best action according to the current Q -values. Now that we are in s_2 after Step 3, what is the probability of seeing the transition (s_2, a_1, s_1) in the next step? That is, calculate the probability of the event “according to the ε -greedy policy, we obtained the action a_1 in the current state s_2 , and after applying this action, the MDP puts us in s_1 as the next state.”
3. If instead of ε -greedy we use the greedy strategy, what action that maximizes Q -values in each state will we choose in the next step?

4 Monte Carlo

Question 8 (2 Points) Consider a two-coin game as shown in the slides for this chapter. In each round, a coin is flipped, and the outcome is either heads or tails. Each round is the same, and the game ends when a coin is heads. You should use the definition of regret over n rounds.

Question 9 (Extra 2 Points) Consider a two-coin game as shown in the slides for this chapter. In each round, a coin is flipped, and the outcome is either heads or tails. Each round is the same, and the game ends when a coin is heads. You should use the definition of regret over n rounds.

In MCTS we use the UCB strategy to balance exploration and exploitation. We expect that the UCB strategy achieves the optimal regret bound for each coin, because the UCB strategy is designed to balance exploration and exploitation.

Give a concrete example of a minimax tree, such that the UCB strategy achieves the optimal regret bound. The MCTS algorithm explores the tree by doing roll-outs from that node.

The mathematical definition of the MCTS algorithm is as follows. The MCTS algorithm does on the tree that you design, so that it is clear that the win rate distribution for the options at the root node may shift over time.

5 Constraint Solving and Propositional Reasoning

Question 10 (3 Points). Consider the following constraint system: The variables are $X = \{x_1, x_2\}$ with $x_1 \in D_1 = [0, 1]$ and $x_2 \in D_2 = [0, 1]$, both are continuous intervals of real numbers. The constraints are

$$C_1 : x_1 = x_2 + 0.2$$

$$C_2 : x_1 = x_2^2$$

Perform propagation on the domain $D = D_1 \times D_2$ using the constraints $\{C_1, C_2\}$ by enforcing arc-consistency. Show how the domains on each variable will be updated in the first several propagation steps, and then show what end results you will obtain by iterating such propagation, and explain why.

145 **Question 11** (2 Points). Put the following formula into CNF:

$$\varphi : \neg \left((p_2 \rightarrow (p_3 \rightarrow \neg p_1)) \rightarrow p_1 \right)$$

146 You only need to expand the definition of the logical connective “ \rightarrow ” and then using De Morgan laws should
147 be enough (i.e., no need to go through the general procedures of introducing new variables, etc.).

