

- You have approximately 170 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 with *circular bubbles*, you should only mark ONE option; for those with *checkboxes*, you should mark ALL that apply (which can range from zero to all options).

| | |
|------------------------------|--|
| First name | |
| Last name | |
| SID | |
| Name of person on your left | |
| Name of person on your right | |

Your Discussion/Exam Prep* TA (fill all that apply):

- | | | | |
|---------------------------------------|---------------------------------------|---------------------------------------|---------------------------------------|
| <input type="checkbox"/> Shizhan (Tu) | <input type="checkbox"/> Peyrin* (Tu) | <input type="checkbox"/> Rachel (W) | <input type="checkbox"/> Mike (W) |
| <input type="checkbox"/> Carl (Tu) | <input type="checkbox"/> Andy (Tu) | <input type="checkbox"/> Henry* (W) | <input type="checkbox"/> Danny* (W) |
| <input type="checkbox"/> Emma (Tu) | <input type="checkbox"/> Wilson (W) | <input type="checkbox"/> Alan (W) | <input type="checkbox"/> Jinkyu (W) |
| <input type="checkbox"/> Mesut* (Tu) | <input type="checkbox"/> Ryan (W) | <input type="checkbox"/> Andreea (W) | <input type="checkbox"/> Lawrence (W) |
| <input type="checkbox"/> Jesse (Tu) | <input type="checkbox"/> Lindsay (W) | <input type="checkbox"/> Chandan (W) | <input type="checkbox"/> Albert (W) |
| <input type="checkbox"/> Cathy (Tu) | <input type="checkbox"/> Gokul* (W) | <input type="checkbox"/> Sherman* (W) | |

| | |
|----------------------------------------------|------|
| Q1. Game Trees | /9 |
| Q2. Short Answer | /14 |
| Q3. Decision Networks and VPI | /9 |
| Q4. Bayes Net CSPs | /9 |
| Q5. Probability and Bayes Net Representation | /20 |
| Q6. Finding Waldo | /12 |
| Q7. Machine Learning: Potpourri | /12 |
| Q8. MDPs and RL | /15 |
| Total | /100 |

To earn the extra credit, one of the following has to hold true. Please circle and sign.

A I spent 170 or more minutes on the practice midterm.

B I spent fewer than 170 minutes on the practice midterm, but I believe I have solved all the questions.

Signature: -----

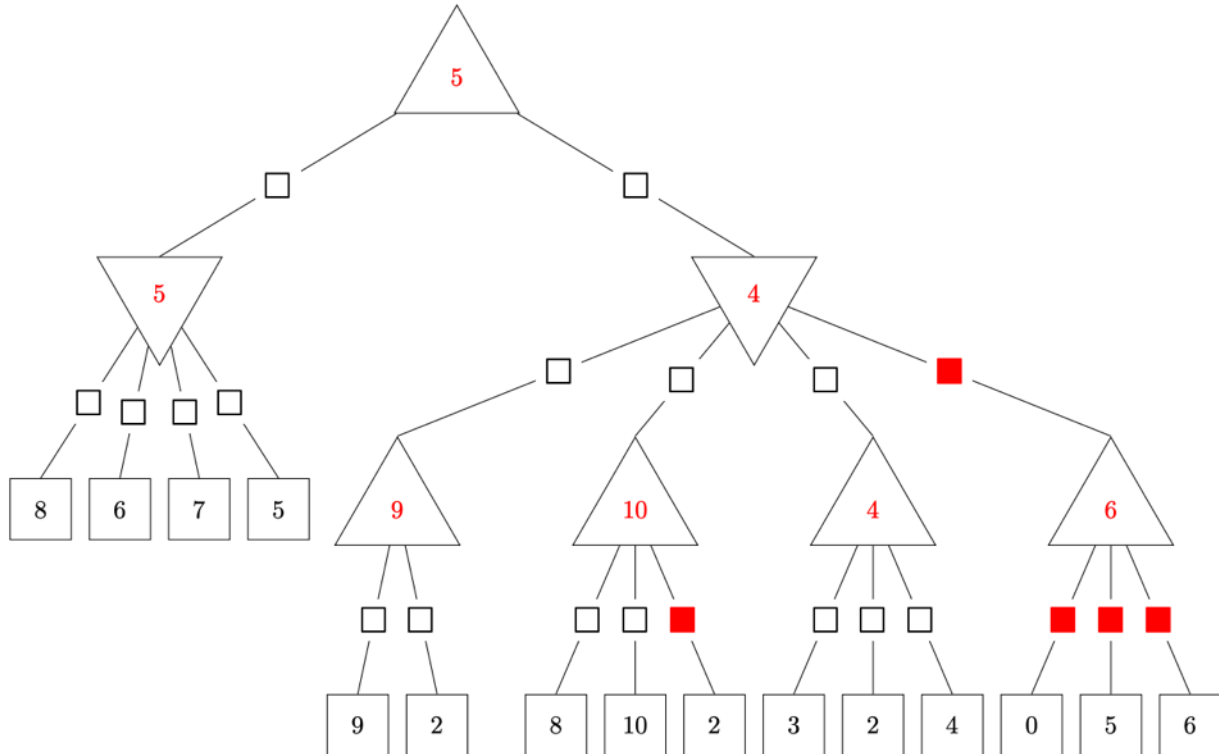
To submit the practice midterm, scan and upload the PDF to Gradescope.

Q1. [9 pts] Game Trees

The following problems are to test your knowledge of Game Trees.

(a) Minimax

The first part is based upon the following tree. Upward triangle nodes are maximizer nodes and downward are minimizers. (small squares on edges will be used to mark pruned nodes in part (ii))



- (i) [1 pt] Complete the game tree shown above by filling in values on the maximizer and minimizer nodes.
- (ii) [3 pts] Indicate which nodes can be pruned by marking the edge above each node that can be pruned (you do not need to mark any edges below pruned nodes). In the case of ties, please prune any nodes that could not affect the root node's value. Fill in the bubble below if no nodes can be pruned.

No nodes can be pruned

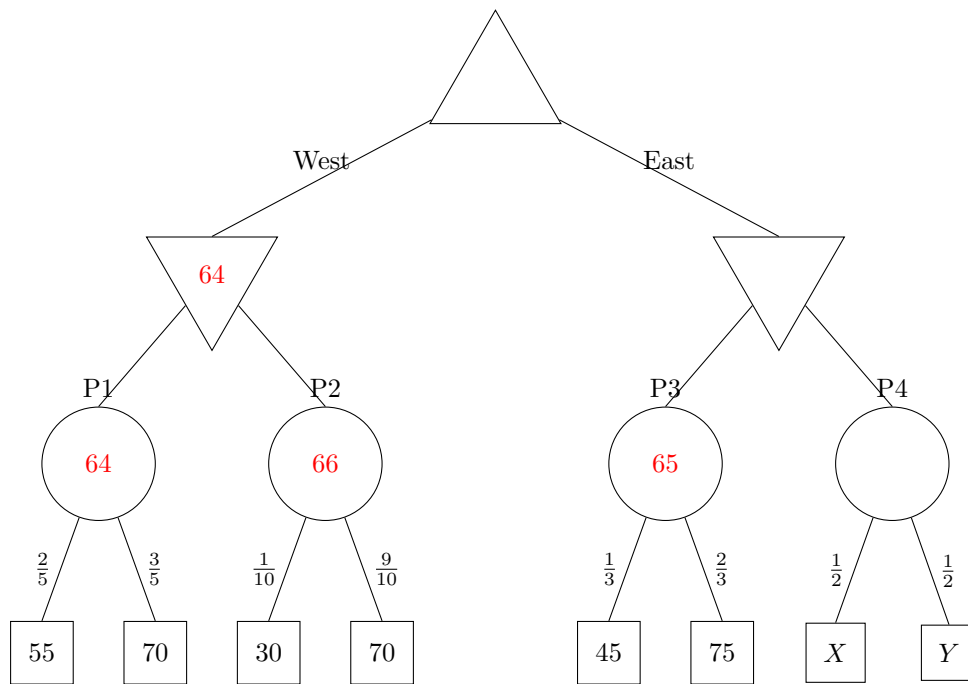
(b) Food Dimensions

The following questions are completely unrelated to the above parts.

Pacman is playing a tricky game. There are 4 portals to food dimensions. But, these portals are guarded by a ghost. Furthermore, neither Pacman nor the ghost know for sure how many pellets are behind each portal, though they know what options and probabilities there are for all but the last portal.

Pacman moves first, either moving West or East. After which, the ghost can block 1 of the portals available.

You have the following gametree. The maximizer node is Pacman. The minimizer nodes are ghosts and the portals are chance nodes with the probabilities indicated on the edges to the food. In the event of a tie, the left action is taken. Assume Pacman and the ghosts play optimally.



(i) [1 pt] Fill in values for the nodes that do not depend on X and Y .

(ii) [4 pts] What conditions must X and Y satisfy for Pacman to move East? What about to definitely reach the P4? Keep in mind that X and Y denote numbers of food pellets and must be **whole numbers**: $X, Y \in \{0, 1, 2, 3, \dots\}$.

To move East: $X + Y > 128$

To reach P4: $X + Y = 129$

The first thing to note is that, to pick A over B , $value(A) > value(B)$.

Also, the expected value of the parent node of X and Y is $\frac{X+Y}{2}$.

$$\implies \min(65, \frac{X+Y}{2}) > 64$$

$$\implies \frac{X+Y}{2} > 64$$

$$\text{So, } X + Y > 128 \implies value(A) > value(B)$$

To ensure reaching X or Y , apart from the above, we also have $\frac{X+Y}{2} < 65$

$$\implies 128 < X + Y < 130$$

$$\text{So, } X, Y \in \mathbb{N} \implies X + Y = 129$$

Q2. [14 pts] Short Answer

- (a) [2 pts] You have a pile of P potatoes to eat and B potato-eating bots. At any time, each bot is either a *chopper* or a *devourer*; all begin as choppers. In a given time step, a *chopper* can *chop*, *idle*, or *transform*. If it chops, it will turn 1 potato into a pile of fries. If it is idle, it will do nothing. If it transforms, it will do nothing that time step but it will be a devourer in the next time step. Devourers are hive-like and can only *devour* or *transform*. When D devourers devour, they will consume exactly D^2 piles of fries that time step – but only if at least that many piles exist. If there are fewer piles, nothing will be devoured. If a devourer transforms, it will do nothing that time step but will be a chopper in the next one. The goal is to have no potatoes or fries left. Describe a minimal state space representation for this search problem. You must write down a size expression in terms of the number of potatoes P , the number of total bots B , the number of fries F , the number of time steps elapsed T , and any other quantities you wish to name. For example, you might write $P^B + T$. You may wish to briefly explain what each factor in your answer represents.

State space size: $P^2 \cdot B$ OR $P \cdot F \cdot B$

P^2 or $P \cdot F$ is need to represent the number of potatoes and fries. It is not sufficient to represent only the number of potatoes remaining, since the devourers can only eat D^2 piles of fries, so number of fries remaining needs to be represented too. B is sufficient to represent the state of the bots because every bot is either a chopper or devourer, their total number is fixed, and only the number of choppers and the number of devourers are relevant to the problem. The individual bot roles do not matter; only the number of each since any chopper can chop and the D devourers devour together.

- (b) [4 pts] Consider a 3D maze, represented as an $(N + 1) \times (N + 1) \times (N + 1)$ cube of $1 \times 1 \times 1$ cells with some cells empty and some cells blocked (i.e. walls). From every cell it is possible to move to any adjacent facing cell (no corner movement). The cells are identified by triples (i, j, k) . The start state is $(0, 0, 0)$ and the goal test is satisfied only by (N, N, N) . Let L_{ij} be the *loose projection* of the cube onto the first two coordinates, where the projected state (i, j) is a wall if (i, j, k) is a wall for *all* k . Let T_{ij} be the *tight projection* of the cube onto the first two coordinates, where the projected state (i, j) is a wall if (i, j, k) is a wall for *any* k . The projections are similarly defined for L_{ik} and so on.

Distance is the maneuver distance. If all paths to the goal are blocked, the distance is $+\infty$.

Mark each admissible heuristic below.

An heuristic h is admissible if it never overestimates the true cost h^* , that is, $h \leq h^*$ for all nodes.

- For (i, j, k) , the value $3N - i - j - k$.

$3(N + 1) - i - j - k$ is the true length of the path from (i, j, k) to the goal. It is the sum of the distance for each dimension, since movement is only possible in straight lines to facing cells.

- For (i, j, k) , the value $N^3 - ijk$.

This is an overestimate. Verify for $N = 4$.

- For (i, j, k) , the distance from (i, j) to the goal in L_{ij} .

The distance to the goal in the projected plane is less than or equal to the distance in 3D. The blocked cells in the plane are those blocked across all of dimension k , and so no further blocks are introduced to lengthen the path.

- For (i, j, k) , the distance from (i, j) to the goal in T_{ij} .

The tight projection blocks a cell in the plane if any cell across dimension k at (i, j) is blocked. In this way it can disconnect the start and goal states, yielding distance $+\infty$, and overestimating the true cost.

- For (i, j, k) , the distance from (i, j) to the goal in L_{ij} plus the distance from (i, k) to the goal in L_{ik} plus the distance from (j, k) to the goal in L_{jk} .

Although the loose projection distance on a particular plane is admissible, the sum of projected distances overcounts the length of the path since the distance in dimensions i and k both appear twice.

- For (i, j, k) , the distance from (i, j) to the goal in T_{ij} plus the distance from (i, k) to the goal in T_{ik} plus the distance from (j, k) to the goal in T_{jk} .

The tight projection distance on a particular plane is not admissible, so the sum is not admissible.

(c) The cube is back! Consider an $(N + 1) \times (N + 1) \times (N + 1)$ gridworld. Luckily, all the cells are empty – there are no walls within the cube. For each cell, there is an action for each adjacent facing open cell (no corner movement), as well as an action *stay*. The actions all move into the corresponding cell with probability p but stay with probability $1 - p$. *Stay* always stays. The reward is always zero except when you enter the goal cell at (N, N, N) , in which case it is 1 and the game then ends. The discount is $0 < \gamma < 1$.

(i) [1 pt] How many iterations k of value iteration will there be before $V_k(0, 0, 0)$ becomes non-zero? If this will never happen, write *never*.

$3N$. At V_0 the value of the goal is correct. At V_1 all cells next to the goal are non-zero, at V_2 all cells next to those are nonzero, and so on.

(ii) [1 pt] If and when $V_k(0, 0, 0)$ first becomes non-zero, what will it become? If this will never happen, write *never*.

$(\gamma p)^{3N} / \gamma$. The value update of a cell c in this problem is $V'(c) = p(r_{c'} + \gamma V(c')) + (1 - p)V(c)$. The first time the value of a state becomes non-zero, the value is $V'(c) = p(r_{c'} + \gamma V(c'))$.

$$V'(g) = p(1 + \gamma V(\text{goal})) = p$$

for a cell g adjacent to the goal.

$$V'(c) = p\gamma V(c')$$

for other cells since the reward is 0.

Carrying out the value recursion, the goal reward +1 is multiplied by p for the step to the goal and $p\gamma$ for each further step. The number of steps from the start to the goal is $3N$, the Manhattan distance in the cube. The first non-zero $V(0, 0, 0)$ is $(\gamma p)^{3N} / \gamma$ since every step multiplies in $p\gamma$ except the last step to the goal, which multiplies in p . Equivalently, the first non-zero value is $p(\gamma p)^{3N-1}$ with $(\gamma p)^{3N-1}$ for all the steps from the start to a cell adjacent to the goal and p for the transition to the goal.

(iii) [1 pt] What is $V^*(0, 0, 0)$? If it is undefined, write *undefined*.

$$\frac{1}{\gamma} \left(\frac{\gamma p}{1 - \gamma + \gamma p} \right)^{3N}$$

To see why, let $V^*(d)$ be the value function of states whose Manhattan distance from the goal is d . By symmetry, all states with the same Manhattan distance from the goal will have the same value. Write the Bellman equations:

$$V^*(d) = \gamma(1 - p)V^*(d) + \gamma pV^*(d - 1) \quad \text{for all } d > 1$$

$$V^*(1) = p + \gamma(1 - p)V^*(1) + \gamma pV^*(0)$$

$$V^*(0) = \gamma V^*(0)$$

and solve starting with $V^*(0) = 0$ to get the answer.

(d) The cube is still here! (It's also still empty.) Now the reward depends on the cell being entered. The goal cell is not special in any way. The reward for *staying* in a cell (either intentionally or through action failure) is always 0. Let V_k be the value function computed after k iterations of the value iteration algorithm. Recall that V_0 is defined to be 0 for all states. For each statement, circle the subset of rewards (if any) for which the statement holds.

(i) [1 pt] As the number of iterations k of value iteration increases, $V_k(s)$ cannot decrease when all cell-entry rewards:

- are zero
 are in the interval $[0, 1]$
 are in the interval $[-1, 1]$
 For zero, the value of every state is constant and zero. For $[0, 1]$, the value can only stay zero or increase. For $[-1, 1]$, this case is identical to $[0, 1]$ except that if the reward to enter a neighboring cell is negative, the *stay* action will be chosen and the value of the cell will stay zero.

(ii) [1 pt] The optimal policy can involve the *stay* action for some states when all cell-entry rewards:

- are zero
 are in the interval $[0, 1]$
 are in the interval $[-1, 1]$
 It is possible for stay to be part of an optimal policy in all three cases. “Can involve” means there exists a set of rewards in the given interval for which there exists an optimal policy that includes the stay action for a state. For all rewards zero, any policy is optimal. For $[0, 1]$ rewards, stay is likewise optimal if the rewards are zero. For $[-1, 1]$ rewards, stay is optimal if rewards are negative (since stay has reward zero) and is optimal if rewards are zero as in the other cases.

(e) F-learning is a forgetful alternative to Q-learning. Where Q-learning tracks Q-values, F-learning tracks F-values. After experiencing an episode (s, a, r, s') , F-learning does the following update:

$$F(s, a) = r + \gamma \max_{a'} F(s', a')$$

As in Q-learning, All F-values are initialized to 0. Assume all states and actions are experienced infinitely often under a fixed, *non-optimal* policy π that suffices for Q-learning's convergence and optimality. Note that π will in general be stochastic in the sense that for each state s , $\pi(s)$ gives a distribution over actions that are then randomly chosen between.

F-learning is equivalent to Q-learning with learning rate $\alpha = 1$.

For each claim, mark the classes of MDPs for which it is true:

(i) [1 pt] F-learning converges to some fixed values:

- for deterministic state transitions
 for stochastic state transitions
 never
 whenever Q-learning converges

(ii) [1 pt] F-learning converges to the optimal Q-values:

- for deterministic state transitions
 for stochastic state transitions
 never
 whenever Q-learning converges

(iii) [1 pt] F-learning converges to the Q-values of the policy π :

- for deterministic state transitions
 for stochastic state transitions
 never
 whenever Q-learning converges

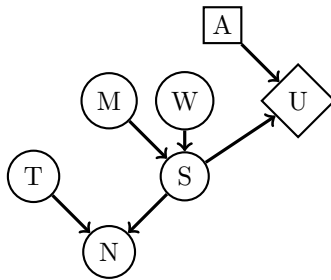
In deterministic MDPs, Q-learning always converges with the usual assumption of all state-actions being experienced infinitely often, and it converges to the optimal values. Learning rate $\alpha = 1$ is actually optimally efficient for the deterministic setting in the sense that Q-learning will converge in the fewest number of steps.

In stochastic MDPs, Q-learning converges when the learning rate is appropriately reduced to 0. F-learning however does not converge: the value is updated to the most recent sample at each step, and so it changes whenever a different transition and reward are experienced.

The policy π is a stochastic policy, so the transitions under this policy. This is true even if the MDP itself is deterministic, since the same action is not necessarily taken in the same state. F-learning never converges in this case since the F-value is updated to the most recent sample each time.

Q3. [9 pts] Decision Networks and VPI

(a) Consider the decision network structure given below:



Mark all of the following statements that **could possibly be true**, for some probability distributions for $P(M), P(W), P(T), P(S|M, W)$, and $P(N|T, S)$ and some utility function $U(S, A)$:

(i) [1.5 pts]

- $VPI(T) < 0$
 $VPI(T) = 0$
 $VPI(T) > 0$
 $VPI(T) = VPI(N)$

VPI can never be negative. $VPI(T) = 0$ must be true since T is independent of S. $VPI(N)$ could also be zero if N and S are independent.

(ii) [1.5 pts]

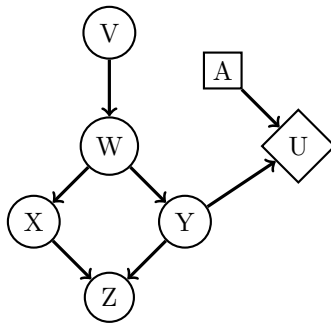
- $VPI(T|N) < 0$
 $VPI(T|N) = 0$
 $VPI(T|N) > 0$
 $VPI(T|N) = VPI(T|S)$

VPI can never be negative. $VPI(T|N) = 0$ if N is conditionally independent of S given N, but will usually be positive. $VPI(T|S) = 0$, and as we've seen $VPI(T|N)$ could also be zero.

(iii) [1.5 pts]

- $VPI(M) > VPI(W)$
 $VPI(M) > VPI(S)$
 $VPI(M) < VPI(S)$
 $VPI(M|S) > VPI(S)$

(b) Consider the decision network structure given below.



Mark all of the following statements that are **guaranteed to be true**, regardless of the probability distributions for any of the chance nodes and regardless of the utility function.

(i) [1.5 pts]

- $VPI(Y) = 0$ *Observing Y could increase MEU*
 $VPI(X) = 0$ *Y can depend on X because of the path through W*
 $VPI(Z) = VPI(W, Z)$ *Consider a case where Y is independent of Z but not independent of W. Then $VPI(Z) = 0 < VPI(W, Z)$*
 $VPI(Y) = VPI(Y, X)$ *After Y is revealed, X will add no more information about Y.*

(ii) [1.5 pts]

- $VPI(X) \leq VPI(W)$ *$VPI(W | X) + VPI(X) = VPI(X, W) = VPI(X | W) + VPI(W)$. We know $VPI(X | W) = 0$, since X is conditionally independent of Y, given W. So $VPI(W | X) + VPI(X) = VPI(W)$. Since VPI is non-negative, $VPI(W | X) \geq 0$, so $VPI(X) \leq VPI(W)$.*

■ $VPI(V) \leq VPI(W)$ Since the only path from V to Y is through W , revealing V cannot give more information about Y than revealing W .

□ $VPI(V | W) = VPI(V)$ $VPI(V | W) = 0$ by conditional independence, but $VPI(V)$ is not necessarily 0

□ $VPI(W | V) = VPI(W)$ Consider a case where W is a deterministic function of V and Y is a deterministic function of W , then $VPI(W | V) = 0 \neq VPI(W)$

(iii) [1.5 pts]

■ $VPI(X | W) = 0$ X is independent of Y given W

□ $VPI(Z | W) = 0$ Y could depend on Z , given W

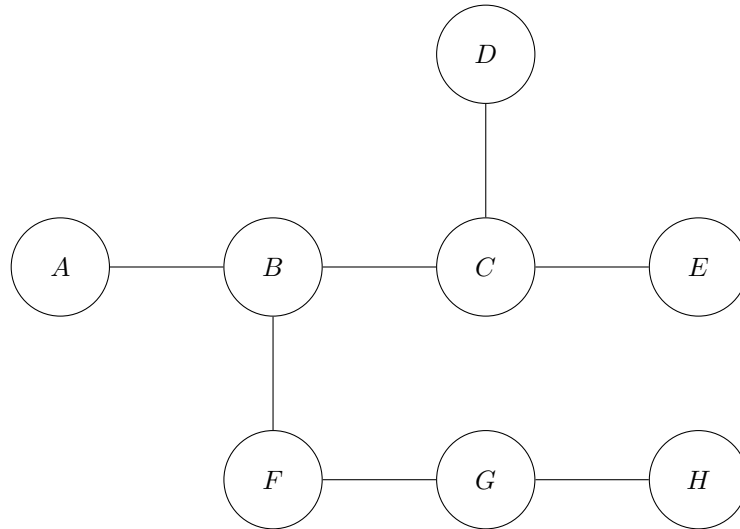
■ $VPI(X, W) = VPI(V, W)$ Both are equal to $VPI(W)$, since both X and V are conditionally independent of Y given W .

□ $VPI(W, Y) = VPI(W) + VPI(Y)$ $VPI(W, Y) = VPI(Y)$, and we can have $VPI(W) > 0$

Q4. [9 pts] Bayes Net CSPs

(a) For the following Bayes' Net structures that are missing a direction on their edges, assign a direction to each edge such that the Bayes' Net structure implies the requested conditional independences and such that the Bayes' Net structure does not imply the conditional independences requested not to be true. Keep in mind that Bayes' Nets cannot have directed cycles.

(i) [2 pts]



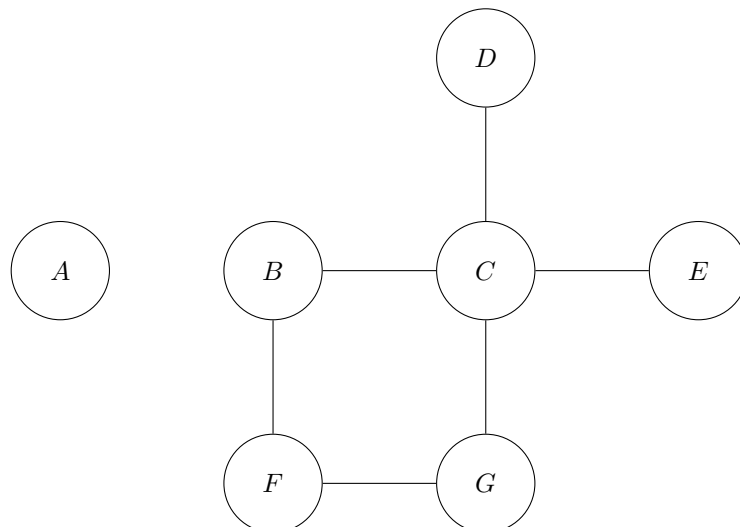
Constraints:

- $D \perp\!\!\!\perp G$
- $D \perp\!\!\!\perp E$
- not $D \perp\!\!\!\perp A$
- $H \perp\!\!\!\perp F$

The following are the directions of the edges:

$B \rightarrow A$
 $C \rightarrow B$
 $D \rightarrow C$
 $E \rightarrow C$
 $F \rightarrow B$
 $F \rightarrow G$
 $H \rightarrow G$

(ii) [2 pts]



Constraints:

- $D \perp\!\!\!\perp F$
- not $D \perp\!\!\!\perp G$

- $D \perp\!\!\!\perp E$
- Bayes Net has no directed cycles

The following are the directions of the edges:

$C \rightarrow B$

$F \rightarrow B$

$F \rightarrow G$

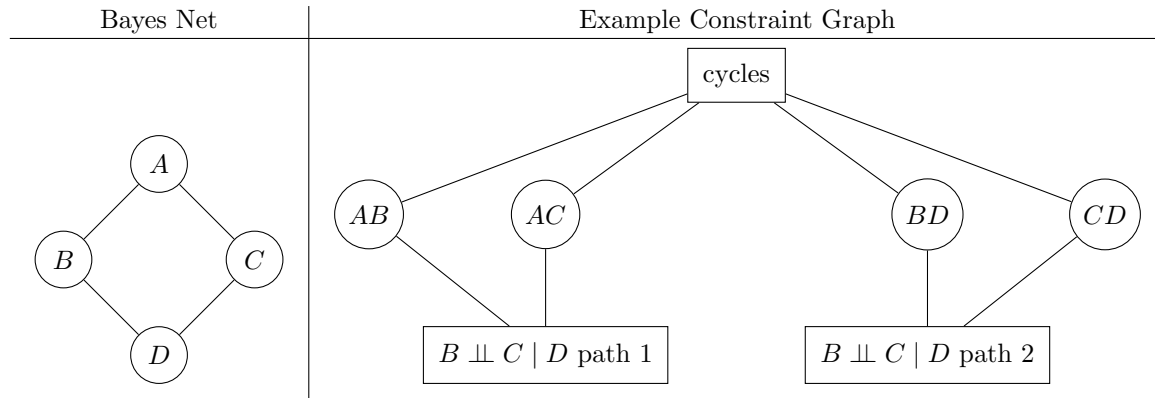
$C \rightarrow G$

$D \rightarrow C$

$E \rightarrow C$

(b) For each of the following Bayes Nets and sets of constraints draw a constraint graph for the CSP. Remember that the constraint graph for a CSP with non-binary constraints, i.e., constraints that involve more than two variables, is drawn as a rectangle with the constraint connected to a node for each variable that participates in that constraint. A simple example is given below.

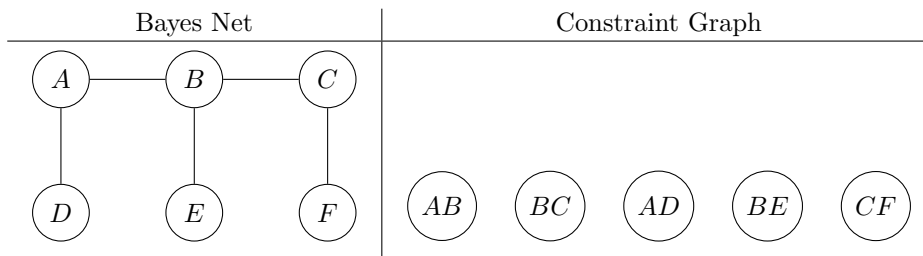
Note: As shown in the example below, if a constraint can be broken up into multiple constraints, do so.



Constraints:

- $B \perp\!\!\!\perp C \mid D$
- No directed cycles

(i) [2 pts]



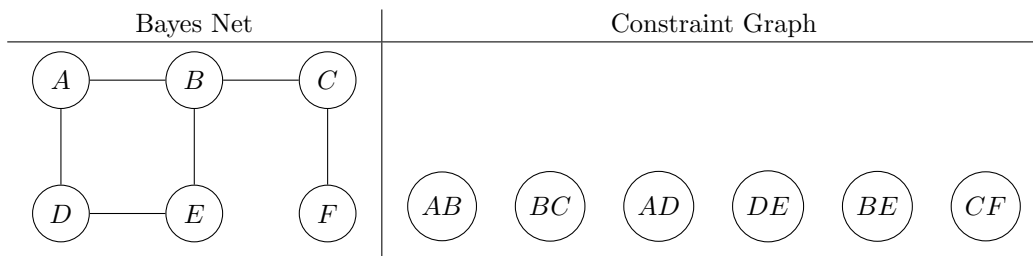
Constraints:

- $A \perp\!\!\!\perp F \mid E$
- not $D \perp\!\!\!\perp C$

Constraint $A \perp\!\!\!\perp F \mid E$: connect AB, BC, BE and CF.

Constraint not $D \perp\!\!\!\perp C$: connect AB, BC and AD.

(ii) [3 pts]



Constraints:

- $A \perp\!\!\!\perp E \mid F$
- $C \perp\!\!\!\perp E$
- No directed cycles

Constraint $A \perp\!\!\!\perp E \mid F$ with path going through path $A - B - E$ with descendant C and F: connect AB, BC, BE, CF.

Constraint $A \perp\!\!\!\perp E \mid F$ with path going through path $A - D - E$: connect AD, DE.

Constraint $C \perp\!\!\!\perp E$ with path going through path $C - B - E$: connect BC, BE.

Constraint $C \perp\!\!\!\perp E$ with path going through path $C - B - A - D - E$: connect AB, BC, AD, DE.

No direct cycles: connect AB, AD, DE and BE.

Q5. [20 pts] Probability and Bayes Net Representation

You're interested in knowing whether you would be **Satisfied** with your choice of snack(s), and so you decide to make the prediction using probabilistic inference over a model with the following variables:

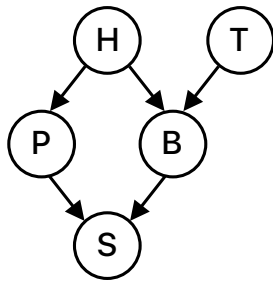
- S , whether or not you will be **Satisfied**.
- H , whether or not you will be **Hungry**.
- T , whether or not you will be **Thirsty**.
- P , whether or not you will have **Pizza**.
- B , whether or not you will have **Boba**.

Each of the variables may take on two values: *yes* or *no*.

- (a) [1 pt] Your first idea for a probability model is a joint probability table over all of the variables. What's the **minimum** number of parameters you need to fully specify this joint probability distribution?

$$2^5 - 1 = 31$$

- (b) [1 pt] You decide this is too many parameters. To fix this, you decide to model the problem with the following Bayes net instead:



| $\Pr(H)$ | |
|----------|-----|
| +h | 0.7 |
| -h | 0.3 |

| $\Pr(T)$ | |
|----------|-----|
| +t | 0.6 |
| -t | 0.4 |

| $\Pr(P H)$ | | |
|------------|----|-----|
| +p | +h | 0.8 |
| +p | -h | 0.5 |
| -p | +h | 0.2 |
| -p | -h | 0.5 |

| $\Pr(B H, T)$ | | | |
|---------------|----|----|-----|
| +b | +h | +t | 0.4 |
| +b | +h | -t | 0.2 |
| +b | -h | +t | 0.9 |
| +b | -h | -t | 0.5 |
| -b | +h | +t | 0.6 |
| -b | +h | -t | 0.8 |
| -b | -h | +t | 0.1 |
| -b | -h | -t | 0.5 |

| $\Pr(S P, B)$ | | | |
|---------------|----|----|-----|
| +s | +p | +b | 0.9 |
| +s | +p | -b | 0.4 |
| +s | -p | +b | 0.7 |
| +s | -p | -b | 0.1 |
| -s | +p | +b | 0.1 |
| -s | +p | -b | 0.6 |
| -s | -p | +b | 0.3 |
| -s | -p | -b | 0.9 |

You do not know which snack(s) you are going for, but you know you are both hungry, thirsty, and definitely getting Pizza. According to your model, what is the probability that you will be satisfied? (First, write out the expression *in terms of conditional probabilities from the model*; then, plug in the *values from the tables* and compute the final answer.)

$$0.6$$

$$\begin{aligned}
 \Pr(+s | +h, +t, +p) &= \sum_b \Pr(+s | +p, b) \cdot \Pr(b | +h, +t) \\
 &= \Pr(+s | +p, +b) \cdot \Pr(+b | +h, +t) + \Pr(+s | +p, -b) \cdot \Pr(-b | +h, +t) \\
 &= 0.9 \times 0.4 + 0.4 \times 0.6 \\
 &= 0.6
 \end{aligned}$$

- (c) [3 pts] You thought the last part required too much computation so you decide to use rejection sampling, sampling variables in topological order. Write the probability of rejecting a sample for the following queries.

$$P(+p \mid +h) =$$

0.3

$$P(-s \mid +p) =$$

$$P(-p \mid +h)P(+h) + P(-p \mid -h)P(-h) = 0.2 \times 0.7 + 0.5 \times 0.3 = 0.29$$

$$P(+s \mid -h, +t) =$$

$$1 - P(-h, +t) = 1 - P(-h)P(+t) = 1 - 0.3 \times 0.6 = 0.82$$

- (d) Given that you are satisfied with your choice of snack(s), write out the variable elimination steps you would take to compute the probability that you actually had boba, that is, $\Pr(+b \mid +s)$. (You do **not** have to plug in the values from the tables.)

- (i) [2 pts] Which of the following factors do we start with?

- $\Pr(H)$ $\Pr(T)$ $\Pr(P)$ $\Pr(B)$ $\Pr(+s)$
 $\Pr(H|P)$ $\Pr(P|H)$ $\Pr(B|H)$ $\Pr(B|T)$ $\Pr(B|H, T)$
 $\Pr(+s|P)$ $\Pr(+s|B)$ $\Pr(+s|P, H)$ $\Pr(+s|P, H, B)$ $\Pr(+s|P, B)$

- (ii) [1 pt] First, we eliminate H . What is the factor f_1 generated when we eliminate H ?

- $f_1(P)$ $f_1(B)$ $f_1(T)$ $f_1(+s)$
 $f_1(P, B)$ $f_1(P, T)$ $f_1(P, +s)$ $f_1(B, T)$ $f_1(B, +s)$ $f_1(T, +s)$
 $f_1(P, B, T)$ $f_1(P, B, +s)$ $f_1(B, T, +s)$

- (iii) [1 pt] Write out the expression for computing f_1 in terms of the remaining factor(s) (before H is eliminated).

$$f_1(P, B, T) = \sum_h \Pr(h) \Pr(P \mid h) \Pr(B \mid h, T)$$

- (iv) [2 pts] Next, we eliminate T . What is the factor f_2 generated when we eliminate T ?

$$f_2(P, B)$$

Write out the expression for computing f_2 in terms of the remaining factor(s) (before T is eliminated).

$$f_2(P, B) = \sum_t \Pr(t) f_1(t, P, B)$$

- (v) [2 pts] Finally, we eliminate P . What is the factor f_3 generated when we eliminate P ?

$$f_3(B, +s)$$

Write out the expression for computing f_3 in terms of the remaining factor(s) (before P is eliminated).

$$f_3(B, +s) = \sum_p \Pr(+s \mid p, B) f_2(p, B)$$

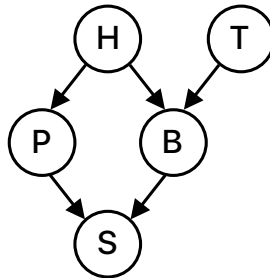
(vi) [1 pt] Write out the expression for computing $\Pr(+b | +s)$ in terms of the remaining factor(s) (after P is eliminated).

$\Pr(+b | +s) =$

$$\frac{f_3(+b, +s)}{\sum_b f_3(b, +s)}$$

(e) **Conditional Independence:** For each of the following statements about conditional independence, mark if it is guaranteed by the Bayes Net.

The Bayes Net is reproduced below for your convenience.



(i) [1 pt] $H \perp\!\!\!\perp T$

- Guaranteed Not guaranteed

(ii) [1 pt] $P \perp\!\!\!\perp T | B$

- Guaranteed Not guaranteed

(iii) [1 pt] $H \perp\!\!\!\perp T | S$

- Guaranteed Not guaranteed

(iv) [1 pt] $S \perp\!\!\!\perp T | B$

- Guaranteed Not guaranteed

(v) [1 pt] $H \perp\!\!\!\perp S | P, B$

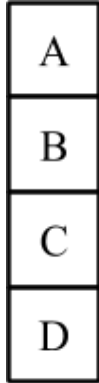
- Guaranteed Not guaranteed

(vi) [1 pt] $P \perp\!\!\!\perp T | H, S$

- Guaranteed Not guaranteed

Q6. [12 pts] Finding Waldo

You are part of the CS 188 Search Team to find Waldo. Waldo randomly moves around floors A, B, C, and D. Waldo's location at time t is X_t . At the end of each timestep, Waldo stays on the same floor with probability 0.5, goes upstairs with probability 0.3, and goes downstairs with probability 0.2. If Waldo is on floor A, he goes down with probability 0.2 and stays put with probability 0.8. If Waldo is on floor D, he goes upstairs with probability 0.3 and stays put with probability 0.7.



| X_0 | $P(X_0)$ |
|-------|----------|
| A | 0.1 |
| B | 0.2 |
| C | 0.3 |
| D | 0.4 |

(a) [2 pts] Fill in the table below with the distribution of Waldo's location at time $t = 1$.

| X_t | $P(X_1)$ |
|-------|--------------------------------------------|
| A | $0.1 * 0.8 + 0.2 * 0.3 = 0.14$ |
| B | $0.2 * 0.5 + 0.1 * 0.2 + 0.3 * 0.3 = 0.21$ |
| C | $0.3 * 0.5 + 0.4 * 0.3 + 0.2 * 0.2 = 0.31$ |
| D | $0.4 * 0.7 + 0.3 * 0.2 = 0.34$ |

(b) [2 pts] $F_T(X)$ is the fraction of timesteps Waldo spends at position X from $t = 0$ to $t = T$. The system of equations to solve for $F_\infty(A)$, $F_\infty(B)$, $F_\infty(C)$, and $F_\infty(D)$ is below. Fill in the blanks.

Note: You may or may not use all equations.

$$\underline{0.8} F_\infty(A) + \underline{0.3} F_\infty(B) + \underline{0} F_\infty(C) + \underline{0} F_\infty(D) = \underline{F_\infty(A)}$$

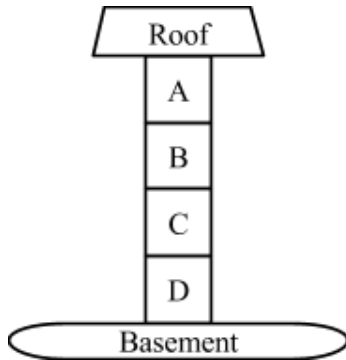
$$\underline{0.2} F_\infty(A) + \underline{0.5} F_\infty(B) + \underline{0.3} F_\infty(C) + \underline{0} F_\infty(D) = \underline{F_\infty(B)}$$

$$\underline{0} F_\infty(A) + \underline{0.2} F_\infty(B) + \underline{0.5} F_\infty(C) + \underline{0.3} F_\infty(D) = \underline{F_\infty(C)}$$

$$\underline{0} F_\infty(A) + \underline{0} F_\infty(B) + \underline{0.2} F_\infty(C) + \underline{0.7} F_\infty(D) = \underline{F_\infty(D)}$$

$$\underline{1} F_\infty(A) + \underline{1} F_\infty(B) + \underline{1} F_\infty(C) + \underline{1} F_\infty(D) = \underline{1}$$

To aid the search a sensor S_r is installed on the roof and a sensor S_b is installed in the basement. Both sensors detect either sound (+s) or no sound (-s). The distribution of sensor measurements is determined by d , the number of floors between Waldo and the sensor. For example, if Waldo is on floor B, then $d_b = 2$ because there are two floors (C and D) between floor B and the basement and $d_r = 1$ because there is one floor (A) between floor B and the roof. The prior of the both sensors' outputs are identical and listed below. **Waldo will not go onto the roof or into the basement.**



| X_0 | $P(X_0)$ |
|-------|----------|
| A | 0.1 |
| B | 0.2 |
| C | 0.3 |
| D | 0.4 |

| S_r | $P(S_r d_r)$ | S_b | $P(S_b d_b)$ |
|-------|-----------------|-------|-----------------|
| +s | $0.3 * d_r$ | +s | $1 - 0.3 * d_b$ |
| -s | $1 - 0.3 * d_r$ | -s | $0.3 * d_b$ |

| S | $P(S)$ |
|-----|--------|
| +s | 0.5 |
| -s | 0.5 |

- (c) [1 pt] You decide to track Waldo by particle filtering with 3 particles. At time $t = 2$, the particles are at positions $X_1 = A$, $X_2 = B$ and $X_3 = C$. Without incorporating any sensory information, what is the probability that the particles will be resampled as $X_1 = B$, $X_2 = B$, and $X_3 = C$, after time elapse?

Answer: $P(X_3 = B|X_2 = A)P(X_3 = B|X_2 = B)P(X_3 = C|X_2 = C)$
 $= (0.2)(0.5)(0.5) = 0.05$

- (d) To decouple this from the previous question, assume the particles after time elapsing are $X_1 = B$, $X_2 = C$, $X_3 = D$, and the sensors observe $S_r = +s$ and $S_b = -s$.

- (i) [3 pts] What are the particle weights given these observations?

| Particle | Weight |
|-----------|-------------------------------------------------|
| $X_1 = B$ | $P(S_r = +s d_r = 1)P(S_b = -s d_b = 2) = 0.18$ |
| $X_2 = C$ | $P(S_r = +s d_r = 2)P(S_b = -s d_b = 1) = 0.18$ |
| $X_3 = D$ | $P(S_r = +s d_r = 3)P(S_b = -s d_b = 0) = 0$ |

- (ii) [1 pt] To decouple this from the previous question, assume the particle weights in the following table. What is the probability the particles will be resampled as $X_1 = B$, $X_2 = B$, and $X_3 = D$?

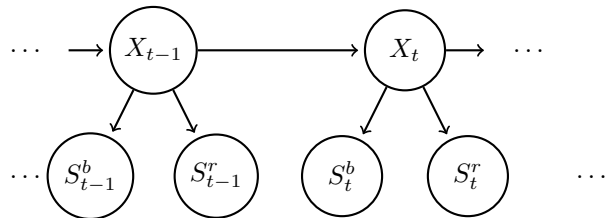
| Particle | Weight |
|----------|--------|
| $X = B$ | 0.1 |
| $X = C$ | 0.6 |
| $X = D$ | 0.3 |

$0.1 * 0.1 * 0.3 = 0.003$

- (e) [3 pts] **Note: the r and b subscripts from before will be written here as superscripts.**

Part of the expression for the forward algorithm update for Hidden Markov Models is given below. $s_{0:t}^r$ are all the measurements from the roof sensor $s_0^r, s_1^r, s_2^r, \dots, s_t^r$. $s_{0:t}^b$ are all the measurements from the roof sensor $s_0^b, s_1^b, s_2^b, \dots, s_t^b$.

Which of the following are correct completions of line (4)? Circle all that apply.



$$P(x_t | s_{0:t}^r, s_{0:t}^b) \propto P(x_t, s_{0:t}^r, s_{0:t}^b) \tag{1}$$

$$= \sum_{x_{t-1}} P(x_{t-1}, x_t, s_{0:t}^r, s_{0:t}^b) \tag{2}$$

$$= \sum_{x_{t-1}} P(x_{t-1}, x_t, s_{0:t-1}^r, s_t^r, s_{0:t-1}^b, s_t^b) \tag{3}$$

$$= \sum_{x_{t-1}} \text{_____} P(x_t | x_{t-1}) P(x_{t-1}, s_{0:t-1}^r, s_{0:t-1}^b) \tag{4}$$

- $P(s_t^r, s_t^b | x_{t-1}, x_t, s_{0:t-1}^r, s_{0:t-1}^b)$
- $P(s_t^r | x_t) P(s_t^b | x_t)$
- $P(s_t^r | x_{t-1}) P(s_t^b | x_{t-1})$
- $P(s_t^r | s_{t-1}^r) P(s_t^b | s_{t-1}^b)$
- $P(s_t^r, s_t^b | x_t)$
- $P(s_t^r, s_t^b | x_t, x_{t-1})$
- None of the above.

There are two equally-correct interpretations of this question: (1) completing the mathematical expression for the probability and (2) completing the algorithmic update for the probability.

Selecting the answers above is correct for interpretation (1): in the Hidden Markov Model, these four probabilities are identical.

Selecting answer 2 alone is correct for interpretation (2): in the Hidden Markov Model, the forward algorithm has the conditional probabilities of the observations given the present state. While the other three choices above are mathematically equivalent, they are not available to the algorithm during execution.

Both correct interpretations earned full credit.

Q7. [12 pts] Machine Learning: Potpourri

- (a) [2 pts] What is the **minimum** number of parameters needed to fully model a joint distribution $P(Y, F_1, F_2, \dots, F_n)$ over label Y and n features F_i ? Assume binary class where each feature can possibly take on k distinct values.

$$2k^n - 1$$

- (b) [2 pts] Under the **Naive Bayes assumption**, what is the **minimum** number of parameters needed to model a joint distribution $P(Y, F_1, F_2, \dots, F_n)$ over label Y and n features F_i ? Assume binary class where each feature can take on k distinct values.

$$2n(k - 1) + 1$$

- (c) [1 pt] You suspect that you are overfitting with your Naive Bayes with Laplace Smoothing. How would you adjust the strength k in Laplace Smoothing?

Increase k

Decrease k

- (d) [2 pts] While using Naive Bayes with Laplace Smoothing, increasing the strength k in Laplace Smoothing can:

Increase training error

Decrease training error

Increase validation error

Decrease validation error

- (e) [1 pt] It is possible for the perceptron algorithm to never terminate on a dataset that is linearly separable in its feature space.

True

False

- (f) [1 pt] If the perceptron algorithm terminates, then it is guaranteed to find a max-margin separating decision boundary.

True

False

- (g) [1 pt] In multiclass perceptron, every weight w_y can be written as a linear combination of the training data feature vectors.

True

False

- (h) [1 pt] For binary class classification, logistic regression produces a linear decision boundary.

True

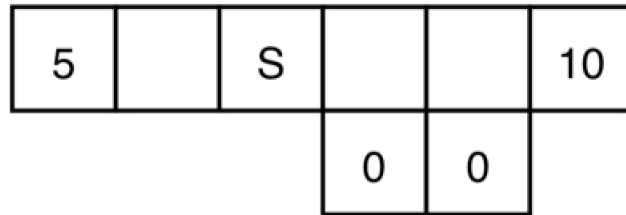
False

- (i) [1 pt] In the binary classification case, logistic regression is exactly equivalent to a single-layer neural network with a sigmoid activation and the cross-entropy loss function.

True

False

Q8. [15 pts] MDPs and RL

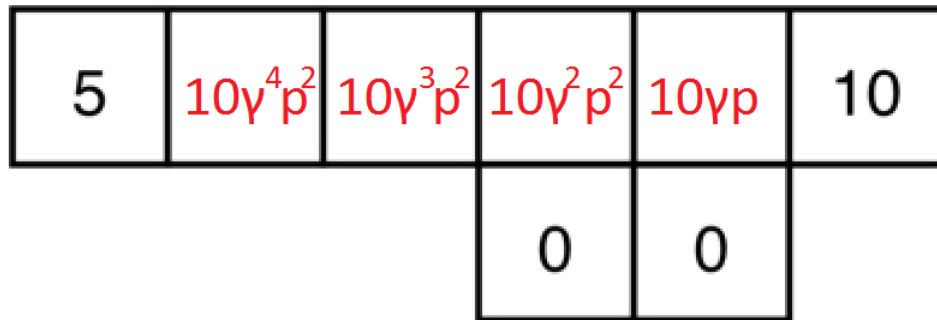


Consider the above gridworld. An agent is currently on grid cell S , and would like to collect the rewards that lie on both sides of it. If the agent is on a numbered square, its only available action is to Exit, and when it exits it gets reward equal to the number on the square. On any other (non-numbered) square, its available actions are to move East and West. Note that North and South are never available actions.

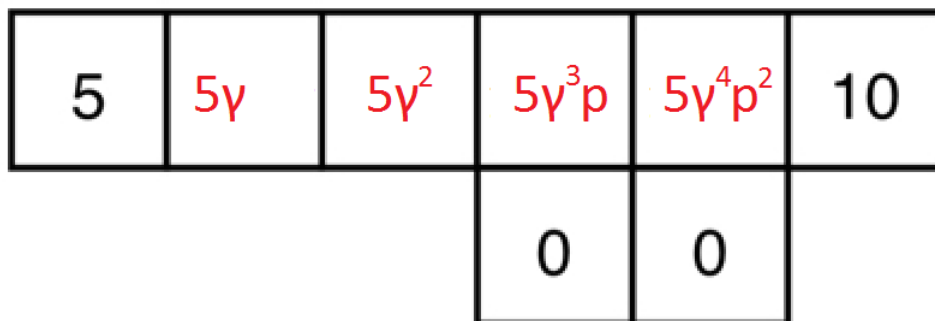
If the agent is in a square with an adjacent square downward, it does not always move successfully: when the agent is in one of these squares and takes a move action, it will only succeed with probability p . With probability $1 - p$, the move action will fail and the agent will instead move downwards. If the agent is not in a square with an adjacent space below, it will always move successfully.

For parts (a) and (b), we are using discount factor $\gamma \in [0, 1]$.

- (a) [2 pts] Consider the policy π_{East} , which is to always move East (right) when possible, and to Exit when that is the only available action. For each non-numbered state x in the diagram below, fill in $V^{\pi_{\text{East}}}(x)$ in terms of γ and p .



- (b) [2 pts] Consider the policy π_{West} , which is to always move West (left) when possible, and to Exit when that is the only available action. For each non-numbered state x in the diagram below, fill in $V^{\pi_{\text{West}}}(x)$ in terms of γ and p .



- (c) [2 pts] For what range of values of p in terms of γ is it optimal for the agent to go West (left) from the start state (S)?

We want $5\gamma^2 \geq 10\gamma^3 p^2$, which we can solve to get:

Range: $p \in [0, \frac{1}{\sqrt{2}\gamma}]$

- (d) [2 pts] For what range of values of p in terms of γ is π_{West} the optimal policy?

We need, for each of the four cells, to have the value of that cell under π_{West} to be at least as large as π_{East} . Intuitively, the farther east we are, the higher the value of moving east, and the lower the value of moving west (since the discount factor penalizes far-away rewards).

Thus, if moving west is the optimal policy, we want to focus our attention on the rightmost cell.

At the rightmost cell, in order for moving west to be optimal, then $V^{\pi_{\text{East}}}(s) \leq V^{\pi_{\text{West}}}(s)$, which is $10\gamma p \leq 5\gamma^4 p^2$, or $p \geq \frac{2}{\gamma^3}$.

However, since γ ranges from 0 to 1, the right side of this expression ranges from 2 to ∞ , which means p (a probability, and thus bounded by 1) has no valid value.

Range: \emptyset

- (e) [2 pts] For what range of values of p in terms of γ is π_{East} the optimal policy?

We follow the same logic as in the previous part. Specifically, we focus on the leftmost cell, where the condition for π_{East} to be the optimal policy is: $10\gamma^4 p^2 \geq 5\gamma$, which simplifies to $p \geq \frac{1}{\sqrt{2}\gamma^3}$. Combined with our bound on any probability being in the range $[0, 1]$, we get:

Range: $p \in \left[\frac{1}{\sqrt{2}\gamma^3}, 1 \right]$, which could be an empty set depending on γ .

Recall that in approximate Q-learning, the Q-value is a weighted sum of features: $Q(s, a) = \sum_i w_i f_i(s, a)$. To derive a weight update equation, we first defined the loss function $L_2 = \frac{1}{2}(y - \sum_k w_k f_k(x))^2$ and found $dL_2/dw_m = -(y - \sum_k w_k f_k(x))f_m(x)$. Our label y in this set up is $r + \gamma \max_a Q(s', a')$. Putting this all together, we derived the gradient descent update rule for w_m as $w_m \leftarrow w_m + \alpha (r + \gamma \max_a Q(s', a') - Q(s, a)) f_m(s, a)$.

In the following question, you will derive the gradient descent update rule for w_m using a different loss function:

$$L_1 = \left| y - \sum_k w_k f_k(x) \right|$$

- (f) [4 pts] Find dL_1/dw_m . Show work to have a chance at receiving partial credit. Ignore the non-differentiable point.

Note that the derivative of $|x|$ is -1 if $x < 0$ and 1 if $x > 0$. So for L_1 , we have:

$$\frac{dL_1}{dw_m} = \begin{cases} -f_m(x) & y - \sum_k w_k f_k(x) > 0 \\ f_m(x) & y - \sum_k w_k f_k(x) < 0 \end{cases}$$

- (g) [1 pt] Write the gradient descent update rule for w_m , using the L_1 loss function.

$$w_m \leftarrow w_m - \alpha dL_1/dw_m$$

$$\leftarrow \begin{cases} w_m + \alpha f_m(x) & y - \sum_k w_k f_k(x) > 0 \\ w_m - \alpha f_m(x) & y - \sum_k w_k f_k(x) < 0 \end{cases}$$