

- You have approximately 170 minutes.
- The exam is open book, open calculator, and open notes.
- For multiple choice questions,
 - ☐ means mark **all options** that apply
 - ☐ means mark a **single choice**

First name	
Last name	
SID	

For staff use only:

Q1.	Tic-Tac-Toe	/11
Q2.	I Want a Project Partner	/20
Q3.	Pac-ML	/23
Q4.	Keyboard Navigation	/15
Q5.	Accommodating Course Robot	/17
Q6.	Holiday Planning	/14
Total		/100

THIS PAGE IS INTENTIONALLY LEFT BLANK

Q1. [11 pts] Tic-Tac-Toe

In this problem we will look at the classic game of Tic-Tac-Toe. Albert is training his intelligent agent, AlbertBot, to compete with MesutBot in a Tic-Tac-Toe competition.

X	O	O
O	X	X
		X

Figure 1: Example winning state for the 'X' player in a Tic-Tac-Toe game

Rules of Tic-Tac-Toe: Two players, X and O, take turns marking an empty square on a three-by-three board. A player wins by placing three marks in a row, horizontally, vertically, or diagonally, as the 'X' player has in the game shown above. If the board fills up with neither player getting three in a row, then the game is a draw.

Assumptions: AlbertBot plays 'O's, and **always plays second**. MesutBot has a *fixed*, unknown, possibly stochastic and possibly suboptimal policy. AlbertBot can train against MesutBot as many times as we want prior to the real competition. In the real competition, AlbertBot gets a reward of 1 if it wins against MesutBot, and a reward of 0 if it draws or loses.

Problem (Goal): Maximize AlbertBot's expected reward, i.e., its expected probability of winning, in the real competition.

This problem is largely motivated from Sutton and Barto's textbook, *Reinforcement Learning: An Introduction*, 2020 new version, section 1.5.

- (a) [2 pts] Among the following choices, which is a tightest upper bound on the number of leaf nodes in the complete game tree of Tic-Tac-Toe, ignoring possible symmetries?

☐ 2^9 ☐ 3^9 ☐ 9^2 ☐ 9^3 ☐ 9^8 ☒ $9!$ ☐ 9^9

The first step has 9 choices, the second step has 8 choices, and so on. Therefore an upper bound is $9 \times 8 \times 7 \times \dots \times 1$

- (b) [2 pts] Is the minimax algorithm a suitable way of solving the problem specified above?

☐ Yes, because the state space of Tic-Tac-Toe is small enough that running Minimax on the whole game tree is feasible.
☐ Yes, but not due to the reason in the choice above.
☐ No, because the minimax only works for zero-sum games, while Tic-Tac-Toe is not zero-sum.
☒ No, but not due to the reason in the choice above.

Minimax finds the optimal policy against the optimal opponent, under which case AlbertBot can never win, so minimax will just return a random policy.

- (c) [3 pts] Let MesutBot's policy be $\pi_M(s, a)$, denoting the probability that MesutBot plays a in state s , and let $Result(s, a)$ denote the state resulting from playing a in s . Can we formulate AlbertBot's Tic-Tac-Toe problem as an MDP, and, if so, which of the following are correct formulas for the transition model $P(s' | a, s)$?

☐ Yes, $P(s' | a, s) = 1$ if $s' = Result(s, a)$, 0 otherwise.
☐ Yes, $P(s' | a, s) = \pi_M(s, a)$ if $s' = Result(s, a)$, 0 otherwise.
☒ Yes, $P(s' | a, s) = \sum_{a': s'=Result(Result(s,a),a')} \pi_M(Result(s, a), a')$.
☒ Yes, $P(s' | a, s) = \pi_M(Result(s, a), a')$ if $s' = Result(Result(s, a), a')$ for some a' , 0 otherwise.
☐ No, AlbertBot's problem cannot be formulated as an MDP.

The last two are equivalent because the rules of Tic-Tac-Toe mean there is at most one action by MesutBot that can cause a transition between any two given states. For a more general case, only the first of the two formulas would be correct.

We now introduce a new method for solving AlbertBot's problem that works as follows:

- (1) Let $V(s)$ be a table, indexed by state, representing the value, i.e., the expected win probability starting in s . Let the initial

values for $V(s)$ be 0.5, except the terminal states: for states with three 'X's in a row, or drawn terminal states, we set an initial value of 0, and for states with three 'O's in a row, we set an initial value of 1.

(2) AlbertBot plays many practice games against MesutBot. In any state s , AlbertBot plays $a^* = \arg \max_a V(\text{Result}(s, a))$, i.e., the best move according to one-step lookahead with V , breaking ties randomly.

(3) For each round (consisting of MesutBot's move followed by AlbertBot's move), let S_t denote the state before the MesutBot's move and S_{t+1} denote the state after AlbertBot's move. V is updated after each round according to the following equation:

$$V(S_t) \leftarrow V(S_t) + \alpha(V(S_{t+1}) - V(S_t))$$

where $V(S_t)$ denotes the value of the state S_t and $0 < \alpha < 1$ is a hyperparameter.

(4) We update V for many games. Then, in the real competition, V remains fixed and AlbertBot again plays $a^* = \arg \max_a V(\text{Result}(s, a))$.

(d) [3 pts] There is at least one major problem in the method as introduced above that can be fixed to make it working better in practice. Which of the following describes the problem, and proposes a reasonable fix to the problem? Choose all correct statements. A statement is considered correct if and only if both the problem and the fix is described correctly.

☐ AlbertBot cannot reason about the opponent (MesutBot's) moves. The fix is to use self-play during training.

☒ AlbertBot cannot reason about the opponent (MesutBot's) moves. The fix is to train a neural network to predict the opponent's move given a certain state.

☐ The training process lacks exploration. The fix is to use epsilon-greedy exploration during training and update the value table for every action taken during training.

☒ The training process lacks exploration. The fix is to use epsilon-greedy exploration during training, but update the value table only for non-exploratory actions during training.

☐ The training process requires full information about the opponent's policy which we do not have access to. The fix is to use Q-learning instead.

☐ The training process requires full information about the opponent's policy which we do not have access to. The fix is to use policy iteration.

☐ None of the above (for every sentence above either the problem or the fix is incorrect).

It is beneficial for Albert to learn a model of the opponent's moving strategy. This can only be learned by playing with the opponent, and not by self-play. The training process does not require full information about opponent's policy. The training process lacks exploration, but we should only update the value table for non-exploratory moves, since in the real competition there will be no exploratory moves.

(e) [1 pt] Given that the problem described in the previous part is successfully addressed, what is a suitable schedule for the value of α ?

☐ Constant close to 1

☐ Constant close to 0

☐ Increase asymptotically towards 1

☒ Decrease asymptotically towards 0

The α value should decrease asymptotically to 0. See the notes on reinforcement learning for details.

Q2. [20 pts] I Want a Project Partner

In 2035, CS 188 will provide an AI partner for your projects. The process of acquiring a partner works as follows:

- In the *Begin* state, you draw a random AI partner whose quality x is an integer drawn from a distribution $P(X)$. $P(X)$ is non-zero over a single contiguous range of values $[a, b]$, where b may be ∞ . (Throughout this question, the variables x and x' will refer to integers in the range $[a, b]$) The cost of this draw is 0.
- Afterwards, you have two choices:
 1. *Stop*: Keep the current partner, and go to the *End* state.
 2. *Draw*: Pay c tokens to draw another AI partner whose quality comes from the same distribution $P(X)$.

We denote state x to be the current quality of your partner. For example, one possible sequence of states and actions might be

$$\text{Begin} \xrightarrow{\text{Draw}} 1 \xrightarrow{\text{Draw}} 2 \xrightarrow{\text{Draw}} 5 \xrightarrow{\text{Stop}} \text{End}$$

with you paying c for each of the second and third draws and getting a quality=5 partner by stopping. The reward for stopping equals the quality of the agent that you get at the end. We assume no discounting, i.e., $\gamma = 1$.

(a) [6 pts] We will go through some steps to solve the problem as an infinite-horizon MDP.

(i) [1 pt] Assume x, x' are positive integers, which of the following are true about the reward function $R(s, a, s')$ of this MDP?

- ☐ For all x , $R(\text{Begin}, \text{Draw}, x) = -c$
☒ For all x, x' , $R(x, \text{Draw}, x') = -c$
☐ For all x , $R(x, \text{Stop}, \text{End}) = -c + x$
☒ For all x , $R(x, \text{Stop}, \text{End}) = x$
☐ None of the Above

From the problem description, you get a free *Draw* in the beginning. After that, each draw has reward $-c$. If you choose to *Stop*, the reward will be the quality of the partner you get.

(ii) [1 pt] What can we say about the Q -values in this MDP?

- ☐ $Q^*(x, \text{Draw})$ increases with x
☒ $Q^*(x, \text{Draw})$ is a constant D , independent of x
☐ $Q^*(x, \text{Draw})$ is a constant D for some (possible empty) sequence $[a, \dots, s]$ and thereafter increases with x
☐ None of the Above

No matter what quality of agent you currently have, the expected return of *Draw* is the same every time.

(iii) [1 pt] True/False If it's optimal to *Stop* at x , it's optimal to *Stop* at any $x' > x$.

- ☒ True
☐ False

If it's optimal to *Stop* at state x , then $U^*(x) = Q^*(x, \text{Stop})$ and $Q^*(x, \text{Stop}) \geq Q^*(x, \text{Draw})$. For $x' > x$, since $Q^*(x', \text{Draw})$ doesn't change, $Q^*(x', \text{Stop}) = x' > x = Q^*(x', \text{Draw})$. Thus, it will be optimal to *Draw* at any x' .

(iv) [1 pt] What can we say about the U -values in this MDP?

- ☒ $U^*(x) = \max_{a \in \{\text{Draw}, \text{Stop}\}} Q^*(x, a)$
☐ $U^*(x)$ increases with x
☐ $U^*(x)$ is a constant D , independent of x
☒ $U^*(x)$ is a constant D for some (possible empty) sequence $[a, \dots, s]$ and thereafter increases with x
☐ None of the Above

First choice is correct because of the definition.

From the explanations in the previous part, we know $U^*(x) = Q^*(x, \text{Stop}) = x$ when x gets very large. We also know that $Q^*(x, \text{Draw})$ is a fixed value. Therefore, for x such that $Q^*(x, \text{Draw}) > x$, it is optimal to *Draw*, and $U^*(x) = Q^*(x, \text{Draw})$, which is a constant.

- (v) [1 pt] Assume at state s , the optimal policy is indifferent between *Draw* and *Stop*, i.e. $Q^*(s, \text{Draw}) = Q^*(s, \text{Stop})$. What can we say about the optimal policy π^* and optimal Q -values Q^* ?

☒ $\pi^*(x) = \text{Draw}$ for all $x < s$

☐ $\pi^*(x) = \text{Stop}$ for all $x < s$

☐ $\pi^*(x) = \text{Draw}$ for all $x > s$

☒ $\pi^*(x) = \text{Stop}$ for all $x > s$

☒ $Q^*(s, \text{Draw}) = s$

☐ None of the Above

Following the previous explanation if at s , we have $Q^*(s, \text{Draw}) = Q^*(s, \text{Stop})$, then it's optimal to *Stop* for $x > s$ and *Draw* for $x < s$. $Q^*(s, \text{Draw}) = s$ because $Q^*(s, \text{Draw}) = Q^*(s, \text{Stop}) = s$.

- (vi) [1 pt] Let $D = Q^*(s, \text{Draw}) = Q^*(s, \text{Stop})$. We are interested in finding D by solving a Bellman equation on U -values. Which, if any, of the following equations are correct?

☐ $D = -c + D \cdot \sum_{x=1}^{\infty} P(x)$

☐ $D = -c + \sum_{x=1}^{\infty} xP(x)$

☐ $D = -c + \sum_{x=1}^D xP(x) + D \cdot \sum_{x=D+1}^{\infty} P(x)$

☒ $D = -c + D \cdot \sum_{x=1}^D P(x) + \sum_{x=D+1}^{\infty} xP(x)$

☐ None of the Above

The bellman equation is $U^*(s) = \max_{a \in \{\text{Draw}, \text{Stop}\}} \sum_x T(s, a, x) \cdot (R(s, a, x) + \gamma U^*(x))$. When $x \in [1, D]$, $U^*(x) = D$. When $x \in [D+1, \infty)$, $U^*(x) = x$.

Therefore, to find D :

$$D = U^*(s) = \sum_{x=1}^D P(x) \cdot (-c + D) + \sum_{x=D+1}^{\infty} P(x) \cdot (-c + x) \\ = -c + D \cdot \sum_{x=1}^D P(x) + \sum_{x=D+1}^{\infty} xP(x)$$

- (b) [4 pts] Now, suppose the $P(X)$ is given in the table below, and suppose $c = \frac{3}{4}$. Remember we define D to be the Q value of the state at which the optimal policy is indifferent between the two actions, i.e. $D = Q^*(s, \text{Draw}) = Q^*(s, \text{Stop})$.

x	$P(x)$
1	1/4
2	1/4
3	1/4
4	1/4

- (i) [2 pts] What is D ? (hint: it's an integer) 2

$$D = -\frac{3}{4} + D \cdot \sum_{x=1}^D P(x) + \sum_{x=D+1}^{\infty} P(x)$$

Plugging value 2, the equation holds.

- (ii) [1 pt] What is $U^*(3)$ 3

Once we figure out $D = 2$, we know $U^*(3) = 3$ because $3 > 2$.

- (iii) [1 pt] What is $\pi^*(1)$ Draw

Once we figure out $D = 2$, we know $\pi^*(1) = \text{Draw}$ because $1 < 2$.

- (c) [3 pts] Now suppose that $c = 1$, for the entire partner-drawing process, $x \in \{1, 2, 3, 4, 5, 6\}$, and the distribution of x , $P(X)$ can be one of $P_{\text{uniform}}(X)$ or $P_{\text{biased}}(X)$. The two distributions are shown in the table below:

x	$P_{\text{uniform}}(x)$	$P_{\text{biased}}(x)$
1	1/6	1/2
2	1/6	1/10
3	1/6	1/10
4	1/6	1/10
5	1/6	1/10
6	1/6	1/10

Define $D_{\text{uniform}} = D$ when $P(X)$ is $P_{\text{uniform}}(X)$, and define $D_{\text{biased}} = D$ when $P(X)$ is $P_{\text{biased}}(X)$. Define Q_{uniform} , Q_{biased} , U_{uniform} , U_{biased} similarly.

(i) [1 pt] Which of the following is true?

- ☒ $D_{\text{uniform}} > D_{\text{biased}}$
☐ $D_{\text{uniform}} < D_{\text{biased}}$
☐ $D_{\text{uniform}} = D_{\text{biased}}$

$D_{\text{uniform}} > D_{\text{biased}}$ because $P_{\text{uniform}}(X)$ has a higher mean, which results in higher expected returns for the action *Draw*.

(ii) [2 pts] Given $D_{\text{biased}} = 2$, which of the following are necessarily false?

- ☒ $U_{\text{uniform}}^*(1) = U_{\text{biased}}^*(1)$
☒ $U_{\text{uniform}}^*(2) = U_{\text{biased}}^*(2)$
☐ $U_{\text{uniform}}^*(3) = U_{\text{biased}}^*(3)$
☐ $U_{\text{uniform}}^*(4) = U_{\text{biased}}^*(4)$
☐ None of the Above

If $D_{\text{biased}} = 2$, and $D_{\text{uniform}} > 2$, then $U_{\text{uniform}}^*(1) > 2$ and $U_{\text{uniform}}^*(2) > 2$. Therefore, they cannot equal to $U_{\text{biased}}^*(1)$ and $U_{\text{biased}}^*(2)$, respectively.

- (d) [7 pts] Suppose that AI agents are produced by two manufacturers, m_1 and m_2 . m_1 manufactures agents whose qualities come from $P_{uniform}(X)$ and m_2 manufactures agents whose qualities come from $P_{biased}(X)$.

Throughout the partner-selecting process, you will be drawing partners from exactly one of the two manufacturers, with an *a priori* probability of 0.5 for each. Moreover, agents manufactured by m_1 and m_2 might have slightly different pitches of voice V , such that agents whose qualities come from $P_{uniform}(X)$ have a 0.8 probability of having “high” pitches while agents whose qualities come from $P_{biased}(X)$ have only a 0.6 probability chance of having “high” pitches.

You are still faced with a decision between *Draw* and *Stop*. However, before you perform any actions, you can listen to the voice of your AI agent. This problem can be formulated as a decision network, as shown in Figure 2.

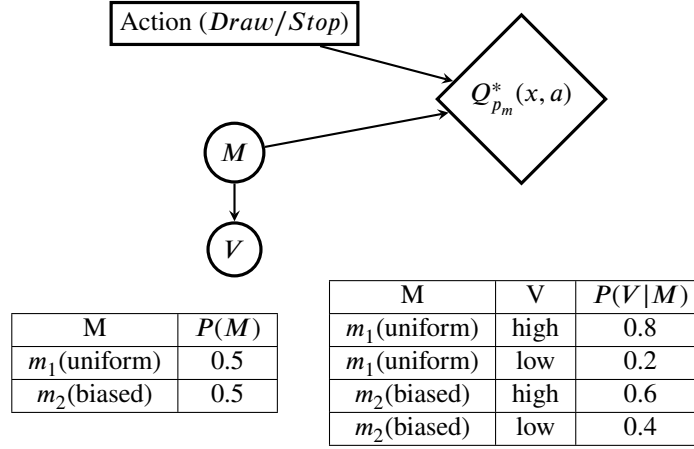


Figure 2: Decision network and associated conditional distributions for the problem.

- (i) [2 pts] Define $EU_x(a|e)$ as “the expected utility of executing action a from state x given evidence e ”, i.e. $EU_x(a|e) = \sum_m P(m|e) \cdot Q_{P_m}^*(x, a)$. Note $Q_{P_m}^*$ represents the optimal Q -values when drawing from the quality distribution P_m associated with manufacturer m , i.e. when $m = m_1$, $Q_{P_m}^*$ is $Q_{uniform}^*$. What is $EU_2(\text{Draw})$?

- ☒ $0.5 \cdot Q_{uniform}^*(2, \text{Draw}) + 0.5 \cdot Q_{biased}^*(2, \text{Draw})$
- ☐ $Q_{uniform}^*(2, \text{Draw})$
- ☐ $Q_{biased}^*(2, \text{Draw})$
- ☐ $\frac{1}{6} \cdot Q_{uniform}^*(2, \text{Draw}) + \frac{1}{10} \cdot Q_{biased}^*(2, \text{Draw})$

$$EU_2(\text{Draw}) = P(M = m_1) \cdot Q_{uniform}^*(2, \text{Draw}) + P(M = m_2) \cdot Q_{biased}^*(2, \text{Draw})$$

- (ii) [2 pts] Define $MEU_x(e)$ as “the maximum expected utility from state x given evidence e ”, i.e. $MEU_x(e) = \max_a EU_x(a|e)$. What is $MEU_2(\emptyset)$?

- ☒ $0.5 \cdot U_{uniform}^*(2) + 0.5 \cdot U_{biased}^*(2)$
- ☐ $U_{uniform}^*(2)$
- ☐ $U_{biased}^*(2)$
- ☐ $\frac{1}{6} \cdot U_{uniform}^*(2) + \frac{1}{10} \cdot U_{biased}^*(2)$

$$\begin{aligned} MEU_2(\emptyset) &= \max_{a \in \{\text{Draw}, \text{Stop}\}} EU_2(a) \\ &= \max_{a \in \{\text{Draw}, \text{Stop}\}} P(M = m_1) \cdot Q_{uniform}^*(2, a) + P(M = m_2) \cdot Q_{biased}^*(2, a) \\ &= P(M = m_1) \cdot U_{uniform}^*(2) + P(M = m_2) \cdot U_{biased}^*(2) \end{aligned}$$

- (iii) [3 pts] Define $VPI_x(E'|e)$ as “the value of observing E' given our current evidence e from state x ”. Given that $D_{uniform} = 3$ and $D_{biased} = 2$, for which state(s) x is $VPI_x(V) > 0$ (Remember V represents the pitches of voice of the AI agents)?

- ☐ 1 ☐ 2 ☐ 3 ☐ 4 ☐ 5 ☐ 6 ☒ None of the above

We start off with $MEU_x(V) = P(V = \text{high}) \cdot MEU_x(V = \text{high}) + P(V = \text{low}) \cdot MEU_x(V = \text{low})$. We know that

1. $MEU_x(V = \text{high}) = P(W = w_1 | V = \text{high}) \cdot U_{uniform}^*(x) + P(W = w_2 | V = \text{high}) \cdot U_{biased}^*(x)$
2. $MEU_x(V = \text{low}) = P(W = w_1 | V = \text{low}) \cdot U_{uniform}^*(x) + P(W = w_2 | V = \text{low}) \cdot U_{biased}^*(x)$

Plugging 1 and 2 into the first equation, we get

$$MEU_x(V) = P(W = w_1, V = \text{high}) \cdot U_{uniform}^*(x) + P(W = w_2, V = \text{high}) \cdot U_{biased}^*(x) + P(W = w_1, V = \text{low}) \cdot U_{uniform}^*(x) + P(W = w_2, V = \text{low}) \cdot U_{biased}^*(x)$$

$$low) * U_{\text{uniform}}^*(x) + P(W = w_2, V = low) * U_{\text{biased}}^*(x)$$

Combining the probabilities, this simplifies to, $P(W = w_1) * U_{\text{uniform}}^*(x) + P(W = w_2) * U_{\text{biased}}^*(x) = \text{MEU}_x(\emptyset)$

Q3. [23 pts] Pac-ML

In this problem we will build classifiers to predict a binary target attribute $y \in \{0, 1\}$.

(a) [4 pts] In this part we will assume that examples are described by a single real-valued attribute x_1 .

(i) [1 pt] Here are three data points:

y	0	0	1
x_1	1	2	3

Are these data points linearly separable?

☒ Yes

☐ No

The separator $x_1 > 2$ classifies the examples correctly.

(ii) [1 pt] Let's add a few more data points in our dataset.

y	0	0	1	0	0	1
x_1	1	2	3	4	5	6

Are these data points linearly separable?

☐ Yes

☒ No

No X_1 -threshold can perfectly separate the data points.

(iii) [2 pts] The data can be augmented by adding a second feature $x_2 = f(x_1)$ for some function f . Which of the following candidates for f will render the data in the previous question linearly separable? For your convenience, the values of $\sin(x)$ when $x \in \{1, 2, 3, 4, 5, 6\}$ are $\{0.84, 0.91, 0.14, -0.76, -0.96, -0.28\}$.

☐ $x_2 = x_1^3$

☐ $x_2 = x_1 \bmod 2$

☐ $x_2 = \sin(x_1)$

☒ $x_2 = x_1 \bmod 3$

Only $x \bmod 3$ can lead to perfect separation. The data points with label $y = 1$ will stay on the x-axis while all the data points labeled $y = 0$ will have a positive x_2 coordinate.

For the remainder of this question we consider examples described by three binary attributes. The training data are as follows:

y	0	0	0	0	0	1	1	1	1	1
x_1	1	0	0	0	0	1	0	0	1	1
x_2	0	0	1	0	0	1	1	1	1	1
x_3	0	0	0	0	0	0	1	1	0	0

(b) [4 pts]

Let us first try a naive Bayes classifier, for which we estimate the necessary probability model \hat{P} from the data.

(i) [1 pt] Given a new data point x'_1, x'_2, x'_3 , which of the following are correct expressions for the predicted label y' of that data point using the naive Bayes classifier? Note that α is a normalizing constant.

☐ $\arg \max_y \alpha \hat{P}(y|x'_1) \hat{P}(y|x'_2) \hat{P}(y|x'_3)$

☒ $\arg \max_y \alpha \hat{P}(y) \prod_{i=1}^3 \hat{P}(x'_i|y)$

☐ $\arg \max_y \alpha \prod_{i=1}^3 \hat{P}(x'_i|y)$

☒ $\arg \max_y \hat{P}(y) \prod_{i=1}^3 \hat{P}(x'_i|y)$

This is just the arg max of the formula for the probability distribution of the label given the evidence using Bayes rule and the "naive" assumption.

(ii) [1 pt] How many parameters need to be estimated for \hat{P} ? Here, "parameters" are the conditional and/or marginal probabilities needed in the naive Bayes model. Do not count parameters that can be calculated using the sum-to-1 constraint.

☐ 4
☒ 7

☐ 8
☐ 10

1 parameter $P(y = 1)$ and two parameters for each $P(x_i|y)$ term, one for the case $y = 1$ and one for $y = 0$.

- (iii) [1 pt] If we discard all the conditional independence assumptions of the naive Bayes model, how many parameters does \hat{P} require? (Again, do not count parameters that can be calculated using the sum-to-1 constraint.)

☐ 9
☐ 12

☒ 15
☐ 16

With no conditional independence, we need the full joint distribution for four variables; $2^4 - 1 = 15$.

- (iv) [1 pt] A new, unlabeled data point arrives with an extra feature x_4 . There are no previous training data to estimate the parameters associated with x_4 , but it is still possible to make a prediction that incorporates x_4 using Laplace smoothing. Predicting y using x_4 and Laplace smoothing in this case will necessarily give identical results to predicting y while ignoring x_4 altogether.

☒ True

☐ False

Laplace smoothing provides pseudocounts for $\hat{P}(x_4|y)$ that are the same for both values of y , so observing x_4 does not affect the posterior for y .

- (c) [3 pts] Now we consider decision trees for classifying these examples. We include the data here as well for convenience.

y	0	0	0	0	0	1	1	1	1	1
x_1	1	0	0	0	0	1	0	0	1	1
x_2	0	0	1	0	0	1	1	1	1	1
x_3	0	0	0	0	0	0	1	1	0	0

- (i) [1 pt] Are all decision trees linear classifiers?

☐ Yes

☒ No

Decision trees can represent any function, including nonlinear ones.

- (ii) [1 pt] Which of x_1 , x_2 , and x_3 has the highest information gain?

☐ x_1

☒ x_2

☐ x_3

Splitting with x_2 first will lead to better split. You don't really need to compute the information gains, you can just visually see that x_2 will lead to an almost perfect split.

- (iii) [1 pt] An adversary would like to flip one of the input bits in the first data column, so that the decision tree learning algorithm cannot fit the training data exactly. Which bit should be flipped?

☐ x_1

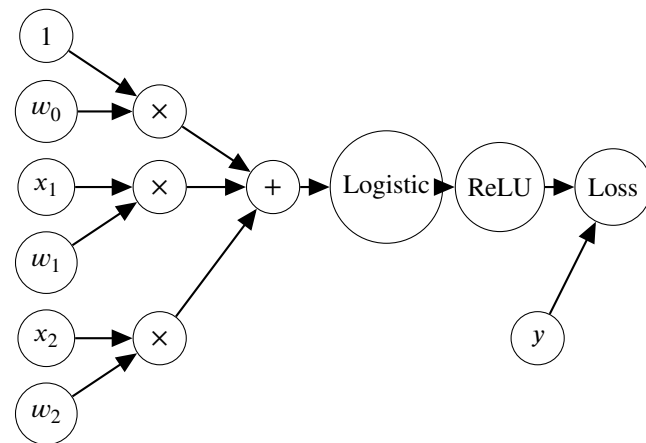
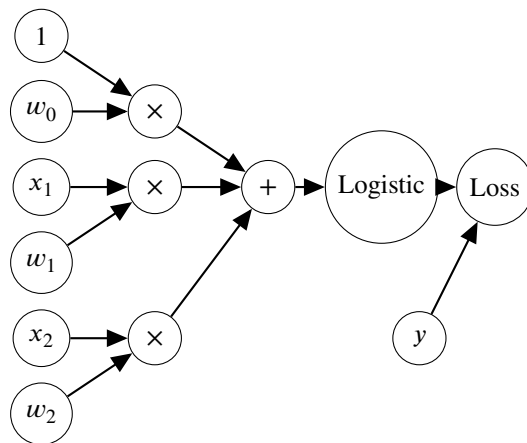
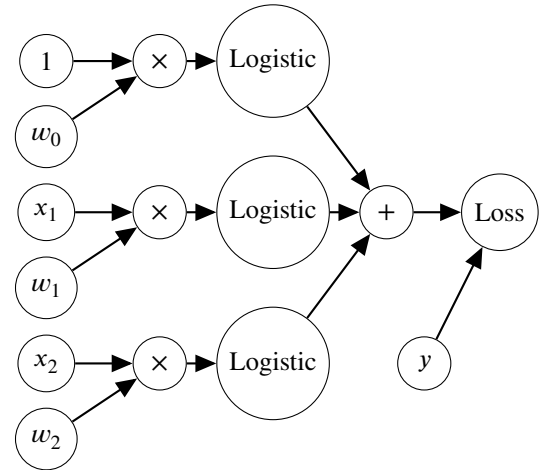
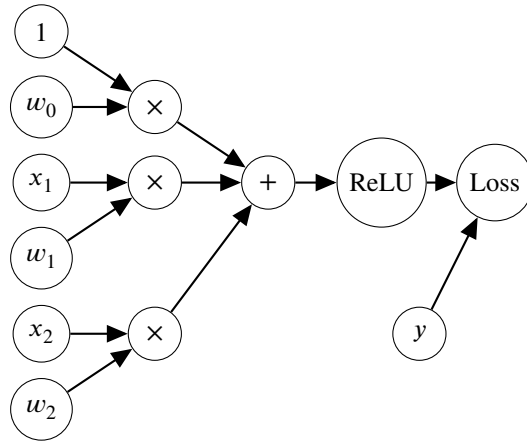
☒ x_2

☐ x_3

The algorithm can always fit the data perfectly *unless* there are two identical examples with different labels. Flipping x_2 to 1 creates such a pair.

- (d) [4 pts] For this part, we consider logistic regression and we ignore the attribute x_3 .

- (i) [2 pts] Which of the following neural networks are capable of representing a logistic regression model, assuming that the "Logistic" nodes implement the logistic function $g(z) = \frac{1}{1+e^{-z}}$?



No ReLUs appear in logistic regression, so the top left is wrong. However, in the bottom right expression the ReLU operation is performed on the logistic function, which is always greater than zero hence $ReLU(logistic(x), 0) = logistic(x)$, i.e. the ReLU ends up being an identity function. The top right choice is wrong as adding the individual logistic functions does not create one function that models a probability of an event.

- (ii) [2 pts] Which of the following expressions is the correct gradient descent update rule for parameter w_2 in the weight vector \mathbf{w} ? Assume that the loss function is the squared loss $\frac{1}{2}(y - h_{\mathbf{w}}(\mathbf{x}))^2$, that $h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$, that $\mathbf{w}^T \mathbf{x}$ denotes the inner-product and that α is the learning rate. Hint: the derivative of the logistic function $g(z) = \frac{1}{1+e^{-z}}$ is $g'(z) = g(z)(1 - g(z))$.

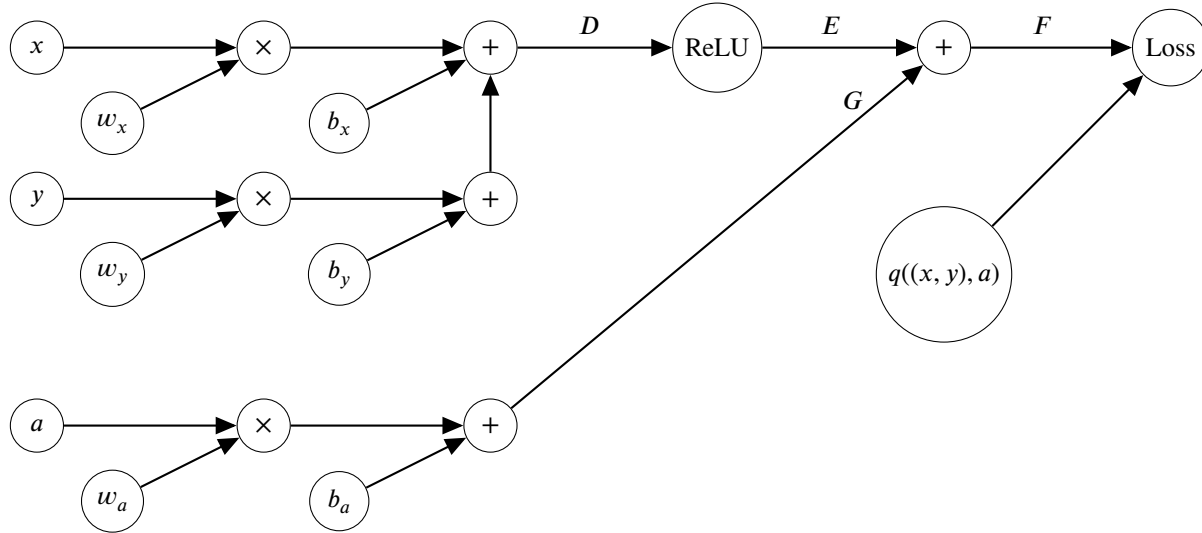
- ☐ $w_2 \leftarrow w_2 - \alpha \left(y - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right) \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \left(1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right) x_2$
☐ $w_2 \leftarrow w_2 - \alpha \left(y - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right) \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \left(y - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right) x_2$
☐ $w_2 \leftarrow w_2 + \alpha \left(y - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right) \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} x_2$
☒ $w_2 \leftarrow w_2 + \alpha \left(y - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right) \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \left(1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right) x_2$

Although we are doing gradient descent, the correct sign is + since we obtain a “−1” from the differentiation of the loss function. The detailed derivation can be found in the textbook.

- (e) In this question we explore using approximate Q-learning to train a robot to complete a fixed, 2D grid maze in the fewest timesteps possible. The Q-function $Q(s, a)$ and policy $\pi(s)$ will be represented by neural networks with input vector $s = [x \ y]$, corresponding to the robot’s x and y coordinates in the grid, and actions $a \in \{1, 2, 3, 4\}$ representing {North, South, East, West} respectively.

Albert trains the neural network shown below to learn his Q-values, where w_x, w_y, w_a are scalar weights and b_x, b_y, b_a are scalar biases.

$q(s = (x, y), a)$ is the target Q-value calculated using transitions (assume that this is the value we want to get close to). The labels D, E, F, G represent the output value of the preceding node:



(i) [2 pts] Which of the following are true about the above neural network? Assume that we use an arbitrary differentiable loss function for the *Loss* node.

- ☐ Using a sigmoid activation instead of the currently used ReLU activation will lead to minimal loss in fewer iterations of training.
- ☐ As long as we train on enough data, we can guarantee a validation accuracy at least as good as our training accuracy.
- ☒ Removing b_x will not affect the expressivity of the neural network.
- ☐ During backpropagation, we compute $\frac{\partial \text{Loss}}{\partial w_a}$, which we then can use to help compute $\frac{\partial \text{Loss}}{\partial w_x}$ and $\frac{\partial \text{Loss}}{\partial w_y}$.
- ☐ None of the above

1. Sigmoid vs. ReLU won't necessarily lead to less or minimal loss.
2. We can't guarantee anything about our validation accuracy regardless of how long we train.
3. b_x and b_y are summed together, and could equivalently be replaced by a single bias term without compromising the expressiveness of the neural network.
4. Although we compute $\frac{\partial \text{Loss}}{\partial w_a}$ in backpropagation, we don't need it to compute $\frac{\partial \text{Loss}}{\partial w_x}$ or $\frac{\partial \text{Loss}}{\partial w_y}$.

(ii) [2 pts] Which of the following are equivalent to $\frac{\partial \text{Loss}}{\partial b_x}$?

- ☒ $\frac{\partial \text{Loss}}{\partial F} \cdot \frac{\partial F}{\partial E} \cdot \frac{\partial E}{\partial D} \cdot \frac{\partial D}{\partial b_x}$
- ☒ $\frac{\partial \text{Loss}}{\partial F} \cdot \frac{\partial F}{\partial E} \cdot \frac{\partial E}{\partial D} \cdot \frac{\partial D}{\partial b_x}$
- ☒ $\frac{\partial \text{Loss}}{\partial F} \cdot \left(\frac{\partial F}{\partial E} \cdot \frac{\partial E}{\partial D} \cdot \frac{\partial D}{\partial b_x} + \frac{\partial F}{\partial G} \cdot \frac{\partial G}{\partial b_a} \cdot \frac{\partial b_a}{\partial b_x} \right)$
- ☒ $\frac{\partial \text{Loss}}{\partial F} \cdot \left(\frac{\partial F}{\partial E} \cdot \frac{\partial E}{\partial D} \cdot \frac{\partial D}{\partial b_x} + \frac{\partial F}{\partial G} \cdot \frac{\partial G}{\partial b_x} \right)$
- ☐ None of the above

All of the above can be simplified into $\frac{\partial \text{Loss}}{\partial F} \cdot \frac{\partial F}{\partial E} \cdot \frac{\partial E}{\partial D}$ because $\frac{\partial D}{\partial b_x} = 1$ and $\frac{\partial b_a}{\partial b_x} = \frac{\partial G}{\partial b_x} = 0$.

(iii) [2 pts] Which of the following loss functions should Albert use for this neural network?

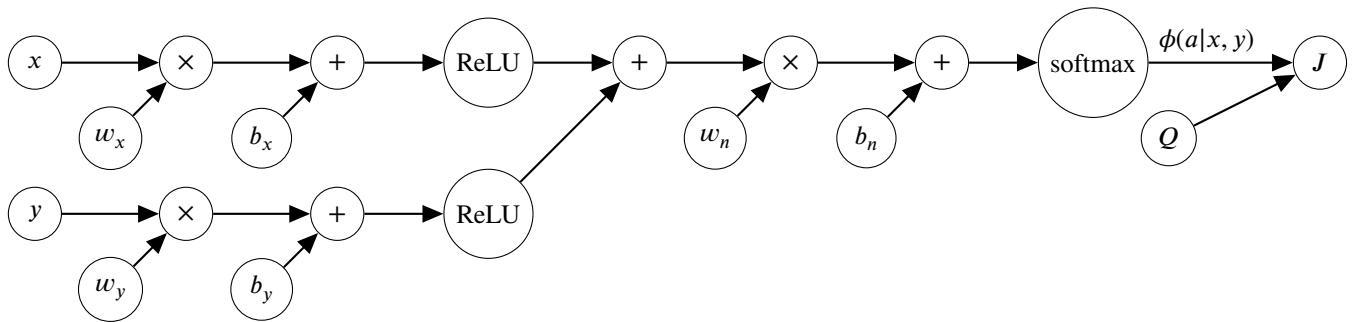
- ☒ Mean-squared-error loss: $L(\hat{y}, y) = (y - \hat{y})^2$
- ☐ $L(\hat{y}, y) = \begin{cases} 0 & \text{if } y = \hat{y} \\ 1 & \text{otherwise} \end{cases}$

- ☐ Cross-entropy loss: $L(\hat{y}, y) = -[t \log y' + (1 - t) \log(1 - y')]$ where $t = \text{zero-one loss}$, $y' = \frac{1}{1+e^{-\hat{y}}}$
- ☐ None of the above

Our goal is to get our neural network to learn how to predict Q-values as close to $q((x, y), a)$.

1. Mean-Squared error loss works for this, because its define as the squared difference between the two values we want to be similar.
2. The indicator function (also known as zero-one loss), won't work because it doesn't give us a reasonable metric for how close our estimates are, and it does not yield a gradient.
3. This is an attempt to apply cross-entropy loss to the problem, which is not suitable since we are not performing classification.

- (f) After training the neural network for $Q(s, a)$, Albert trains another neural network to return a distribution of actions to take at the current state $\phi(a|s)$. w_n is a vector weight in \mathbb{R}^4 and b_n is a vector bias in \mathbb{R}^4 . This neural network is shown below:



The node Q in the above represents the neural network from the previous part, which takes in (x, y, a) and returns a Q-value. $\phi(a|x, y)$ is a vector containing a probability distribution over all possible actions the robot can take from a state. Albert defines J as a function that computes the expected Q-value weighted by the probabilities of the policy network output as follows:

$$J(\phi, Q, x, y) = \sum_a \phi(a|x, y) Q((x, y), a)$$

We want to train $\phi(a|x, y)$ to put higher probability on actions associated with larger Q-values.

- (i) [1 pt] Which of the following can we use to appropriately update w_y at every iteration?

- ☐ Gradient descent
- ☒ Gradient ascent
- ☐ Stochastic gradient descent
- ☐ Perceptron learning rule
- ☐ None of the above

We're maximizing a likelihood function, so we use gradient ascent.

- (ii) [1 pt] Which of the following is / are the appropriate update rule(s) for w_x based on the above neural network and J function? We set α to be the learning rate.

- ☐ $w_x \leftarrow w_x - \alpha \nabla_{w_x} J$
- ☒ $w_x \leftarrow w_x + \alpha \nabla_{w_x} J$
- ☐ $w_x \leftarrow w_x + \alpha (\nabla_{w_x} J + \nabla_{b_x} J)$
- ☐ $w_x \leftarrow w_x - \alpha (\nabla_{w_x} J - |\nabla_{b_x} J|)$
- ☐ None of the above

This follows from the gradient ascent update rule, where we sum the previous weight with the gradient of the likelihood function times the learning rate.

Q4. [15 pts] Keyboard Navigation

Pacman is navigating in a map shown in the figure below (following the QWERTY keyboard). There are $N = 12$ grid squares in total. Each square is labeled with a letter and its coordinate. Let $C_t = (x_t, y_t)$ denote Pacman's position at time step t , and it can move to any of the 6 directions to $(x_t, y_t + 1)$, $(x_t, y_t - 1)$, $(x_t - 1, y_t)$, $(x_t + 1, y_t)$, $(x_t + 1, y_t + 1)$, or $(x_t - 1, y_t - 1)$ in the next timestep. When Pacman attempts to move out of the map, i.e., $x_{t+1} \notin \{1, 2, \dots, 6\}$ or $y_{t+1} \notin \{1, 2\}$, it will stay in the same square.

We assume that Pacman knows the layout of the map. However, Pacman does not know its position, and also does not know whether the move succeeds (i.e., not staying in the same square) at each step. There are two sensors to help Pacman track its position: (1) A_t is the number of neighboring squares for (x_t, y_t) . For example, $A_t = 4$ when $(x_t, y_t) = (2, 2)$. (2) E_t is the number of alphabetically adjacent letters among the neighboring squares. The table below summarizes the (deterministic) value of E_t for each square.

Q(1,2)	W(2,2)	E (3,2)	R (4,2)	T (5,2)	Y(6,2)
A(1,1)	S(2,1)	D (3,1)	F (4,1)	G(5,1)	H (6,1)

C_t	E_t
Q, W, R, Y, A, S, T	0
E, D, F, H	1
G	2

- (a) At step $t = 0$, Pacman can be in any square of the map with equal probability. In this part, at each step t , Pacman can move in any of the 6 directions with an equal probability. Note that Pacman can stay in the same square by taking an invalid move out of the map. We denote $Pr[C_t = (x_t, y_t)]$ as the probability of Pacman at position (x_t, y_t) at step t .

- (i) [1 pt] Assume that Pacman does not have access to sensor signals. What is $Pr[C_1 = (2, 2)]$, i.e., the probability that Pacman is at position (2, 2) when $t = 1$? Please fill in your solution as a fraction below. Please reduce the fraction into the lowest terms. For example, $2/72$ should be simplified to $1/36$. Please fill in 0/1 if the answer is 0.

/

Probability of Start at Q and move to W: $1/12 * 1/6 = 1/72$. Similarly, probability of start at A or S or E and move to W is also $1/72$. The probability of starting at W and remaining at W due to failed action is $1/12 * 2/6 = 1/36$. The total probability is $1/72 * 4 + 1/36 = 1/12$

- (ii) [1 pt] Pacman observes that $E_0 = 1$. Given this evidence, what is $Pr[C_1 = (2, 2) | E_0 = 1]$?

/

Because $E_0 = 1$, $C_0 = E, D, F, H$ with probability $1/4$ on each. Only if $C_0 = E$, taking one of the six actions takes Pacman from E to W. So the probability of Pacman "at S at time step 1" is $1/4 * 1/6 = 1/24$.

- (b) Pacman decides to track its location by particle filtering. Let $(x_t^{(i)}, y_t^{(i)})$ be the location of the i -th particle at step t . At step $t = 0$, we have 4 particles initialized in $\{(1, 1), (1, 2), (6, 1), (6, 2)\}$.

- (i) [2 pts] Assume that after one forward simulation update, these particles reach states $\{(2, 1), (2, 2), (6, 2), (6, 2)\}$ respectively. What is the probability of each of these transitions?

Particle 1 from (1, 1) to (2, 1): /

Particle 4 from (6, 2) to (6, 2): /

The transition probabilities from (1, 1) to (2, 1) is $1/6$. The transition probabilities from (6, 2) to (6, 2) (probability of failed action) is $3/6$ or $1/2$.

- (ii) [2 pts] What is the updated belief distribution based on the new particle values? Note that we have not incorporated any observations.

$Pr[C_1 = (2, 1)] =$ /

$Pr[C_1 = (2, 2)] =$ /

$$Pr[C_1 = (6, 2)] = \boxed{1} / \boxed{2}$$

Among the 4 particles, 1 is at (2, 1), 1 is at (2, 2), 2 are at (6, 2), which results in the answer.

- (iii) [2 pts] Pacman observes that $A_1 = 4$. Please fill in the probabilities given the evidence.

$$Pr[C_1 = (2, 1) | A_1 = 4] = \boxed{1} / \boxed{2}$$

$$Pr[C_1 = (2, 2) | A_1 = 4] = \boxed{1} / \boxed{2}$$

$$Pr[C_1 = (6, 2) | A_1 = 4] = \boxed{0} / \boxed{1}$$

The 2 particles at (6, 2) disagrees with the observation, thus have a weight of 0. The other two particles have a weight of 1.

- (c) In this part, Pacman's goal is to reach the position $(g_x, g_y) = (1, 2)$ with the minimal cost. The cost of a move from (x_t, y_t) to (x_{t+1}, y_{t+1}) is $1 + E_{t+1}$. For example, a move from (3, 2) to (2, 2) has the cost of 1, and a move from (2, 2) to (3, 2) has the cost of 2. Note that even if Pacman stays in the same square after the move, Pacman still needs to pay the cost; e.g., moving left from (1, 1) has the cost of 1. Knowing its position at step $t = 0$, Pacman performs a search to find the optimal sequence of actions to take. Let $h(x, y)$ be the cost from (x, y) to (g_x, g_y) .

- (i) [3 pts] Which of the following are admissible heuristics for $h(x, y)$? Note that we know $(g_x, g_y) = (1, 2)$, i.e. the letter Q.

☒ $|x - g_x| + |g_y - y|$

☒ $\min(|x - g_x|, |g_y - y|)$

☒ $\max(|x - g_x|, |g_y - y|)$

☒ $2 \min(|x - g_x|, |g_y - y|)$

☐ $\min(|x - g_x|, |g_y - y|) + 1$

☐ $\max(|x - g_x|, |g_y - y|) + 1$

☒ $\max(|x - g_x|, |g_y - y|) + \mathbf{I}[x > 3]$, where $\mathbf{I}[x > 3] = 1$ when $x > 3$, and $\mathbf{I}[x > 3] = 0$ otherwise.

☐ None of the above

The cost is always greater than or equal to the Manhattan distance, to admissible heuristics to the Manhattan distance works. The last choice also works because both E and D has cost value 2.

- (ii) [2 pts] In this part and onwards, Pacman does not know the position at step $t = 0$, but can infer possible positions according to the evidence.
Pacman observes that $A_0 = 2$ and $E_0 = 1$ at step $t = 0$. What is the cost of the path found by A* Tree Search with an admissible heuristics?

$$\boxed{7}$$

By the observations we can directly know that the starting position is H. An optimal path from H to Q is HYTREWQ which has cost 7.

- (iii) [2 pts] Pacman observes that $A_0 = 3$ and $E_0 = 0$ at step $t = 0$. We want to find an action sequence such that Pacman is guaranteed to be in $(g'_x, g'_y) = (1, 1)$ after performing this action sequence. We also want to minimize the worst-case cost of this sequence among all possible initial positions (x_0, y_0) . Pacman does not have access to A_t or E_t when $t > 0$. What is the worst-case cost of this action sequence among all possible initial positions?

$$\boxed{6}$$

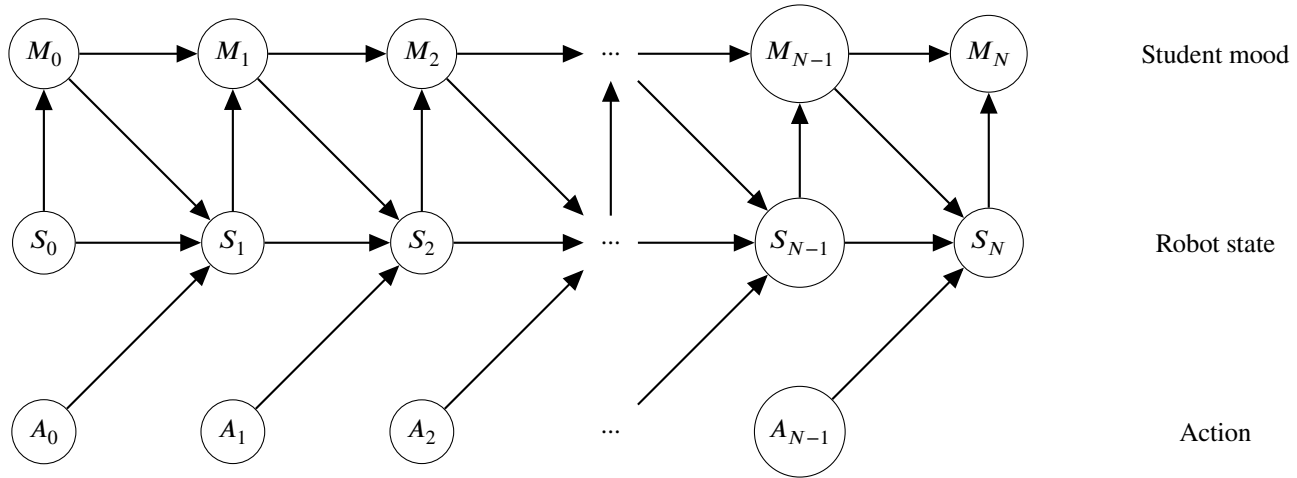
By the observations, we know that the starting position is A or Y. Starting from A would give us a cost of 0, since we are already at the goal. If we start at Y we can follow the path YTREWA with cost 6. The same sequence of actions can make us stay at A if we initially start at A. Thus cost 6 is the worst case cost that guarantees we end up at A = (1, 1).

Q5. [17 pts] Accommodating Course Robot

The CS188 GSIs are designing an instructional robot to improve the mood of the students.

At every timestep t , the robot stores in its state a list of ongoing CS188 assignments S_t , and the robot takes action A_t to release a list of future assignments S_{t+1} . The released assignments S_{t+1} are also improved by student mood and feedback from the previous timestep t , denoted M_t .

This gives us the Bayes net below.



- (a) [1 pt] Which of the below represent the conditional probability table for the variable S_{t+1} in the Bayes net shown above?

We denote this quantity $T(S_{t+1})$, as it is the transition model for S at time $t + 1$.

- ☐ $P(S_{t+1}|S_t, A_t)$
☒ $P(S_{t+1}|S_t, A_t, M_t)$
☐ $P(S_{t+1}|S_t, A_t, M_{t-1})$
☐ None of the above

Note that the Bayes net conditional probability tables are $P(X|\text{parents of } X)$. In this case, the parents of S_{t+1} are S_t, M_t, A_t . This gives us $P(S_{t+1}|S_t, A_t, M_t)$.

- (b) [1 pt] Which of the below represent the conditional probability table for the variable M_{t+1} with the Bayes net shown above? We denote this quantity $T(M_{t+1})$, as it is the transition model for M at time $t + 1$.

- ☐ $P(M_{t+1}|M_t)$
☐ $P(M_{t+1}|M_t, A_t)$
☐ $P(M_{t+1}|S_t, A_t)$
☐ $P(M_{t+1}|S_t)$
☐ $P(M_{t+1}|S_{t+1})$
☒ $P(M_{t+1}|M_t, S_{t+1})$
☐ None of the above

The only parents of M_{t+1} are M_t, S_{t+1} . Hence the CPT is $P(M_{t+1}|M_t, S_{t+1})$.

- (c) [2 pts] What is the joint transition model for variables (S_t, M_t) in the Bayes net above? We denote this quantity $T(S_{t+1}, M_{t+1})$. Include all choices that can be proved to be equivalent using conditional independence assertions expressed by the Bayes net.

- ☒ $P(S_{t+1}, M_{t+1}|S_t, M_t, A_t)$
☐ $P(S_{t+1}, M_{t+1}|S_t, M_t)$
☒ $T(S_{t+1})T(M_{t+1})$
☒ $P(S_{t+1}, M_{t+1}|S_{1:t}, M_{1:t}, A_{1:t})$
☐ $P(S_{t+1}, M_{t+1}|S_{1:t})$

- ☐ $P(S_{t+1}, M_{t+1} | A_{1:t})$
☐ $P(S_{t+1}, M_{t+1} | M_{1:t})$
☐ None of the above

The transition model is defined to be:

$$\begin{aligned}
 P(S_{t+1}, M_{t+1} | S_t, M_t, A_t) &= P(S_{t+1} | S_t, M_t, A_t) P(M_{t+1} | S_{t+1}, S_t, M_t, A_t) \text{ (product rule)} \\
 &= P(S_{t+1} | S_t, M_t, A_t) P(M_{t+1} | M_t, S_{t+1}) \text{ (conditional independence)} \\
 &= T(S_{t+1}) T(M_{t+1})
 \end{aligned}$$

By conditional independence, we have $P(S_{t+1}, M_{t+1} | S_t, M_t, A_t) = P(S_{t+1}, M_{t+1} | S_{1:t}, M_{1:t}, A_{1:t})$

(d) [2 pts] Is the above Bayes net Markovian in (S_t, M_t) ?

- ☒ Yes, always Markovian
☐ No, never Markovian
☐ It is possibly Markovian, depending on whether the values $m_{1:t}$ are given
☐ None of the above

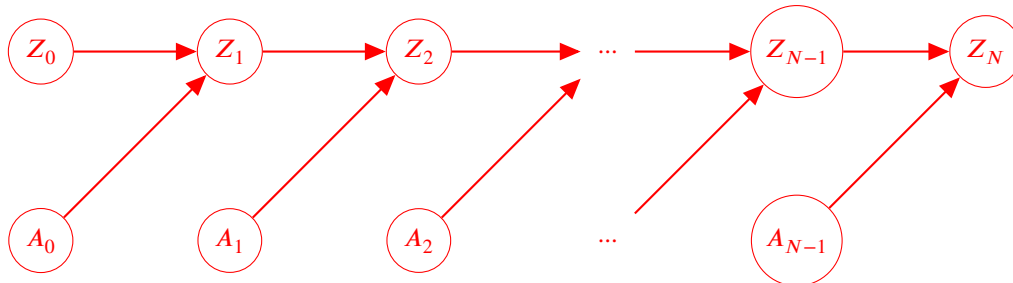
We need to check if past variables are independent of future variables given the present variables. This is indeed true, since $\forall k > 0, (S_{t+k}, M_{t+k})$ is conditionally independent of $S_{0:t-1}$ and $M_{0:t-1}$ given (S_t, M_t) .

(e) [2 pts] You decide the above Bayes net is too complex, and want to try to simplify it down to the MDP format we've learned in class, where the transition dynamics are given by $P(z_{t+1} | z_t, a_t)$ for some single state variable z_t , instead of the transition functions you computed in earlier subparts of this problem. What are valid approaches to do this?

- ☐ No modifications needed.
☒ Redefine states to be $z_t = [s_t, m_t]$
☒ Redefine states to be $z_t = [s_t, m_{t-1}]$, $\forall t \geq 1$
☒ Perform variable elimination on m_t , $\forall t$
☒ Perform variable elimination on s_t , $\forall t$
☐ None of the above

To use the transition function $T(s_t, a_t, s_{t+1})$, we need to make the next state Markovian in the sense that it depends only on the previous state and previous action. All choices that merge the s_t and m_t (as long as they are connected by an edge) into a single state tuple, will give us the desired behavior.

Note that if we do variable elimination on the m_t or the s_t 's, this equivalently gives us a Bayes net of the following form. $z_t = m_t$ or s_t , depending on which one was left after eliminating.



This Bayes net exactly models the MDP seen during class, where the transition function is $P(z_{t+1} | z_t, a_t)$

For the remainder of the question, we denote the transitions in this setting as $(s_t, m_t, a_t, r_t, s_{t+1}, m_{t+1})$, where $r_t = R(m_{t+1})$, since the reward is entirely dictated by the mood of the student, and not the state of the robot. We use discount factor γ .

(f) [2 pts] What is the discounted sum of rewards if the robot observes a sequence of student moods m_0, \dots, m_N ?

- ☒ $\sum_{t=0}^{N-1} \gamma^t R(m_{t+1})$
☐ $R(m_1) + \gamma V([s_2, m_2])$
☐ $R(m_1) + \gamma \max_{a_1} Q([s_2, m_2], a_1)$
☐ None of the above

At timestep 0, the robot receives reward $R(m_1)$. At timestep 1, the discounted reward received is $\gamma^1 R(m_2)$. The very last

timestep, the discounted reward term is $\gamma^{N-1} R(m_N)$. Hence the discounted sum of rewards is $\sum_{t=0}^{N-1} \gamma^t R(m_{t+1})$

Note that the choices involving the V and Q functions are not correct as those assume acting optimally after the 0th timestep, and acting optimally is not guaranteed to give us the sequence of student moods observed, m_0, \dots, m_N .

- (g) [2 pts] Say you wish to run Q-learning on some dataset of transitions $D = \{(s_t, m_t, a_t, r_t, s_{t+1}, m_{t+1})\}_{t=0}^{N-1}$. Each transition is generated by picking a state, taking an action, and observing the result.

You initialize all $Q([s_t, m_t], a_t) = 0, \forall (s_t, m_t), a_t$, and then use the below approach, largely plugging in the Q-learning algorithm from lecture:

1. Randomly select some transition from the dataset $(s_t, m_t, a_t, r_t, s_{t+1}, m_{t+1})$.
2. Compute $sample = r_t + \gamma \max_{a_{t+1}} Q([s_{t+1}, m_{t+1}], a_{t+1})$
3. $Q([s_t, m_t], a_t) \leftarrow \alpha(sample) + (1 - \alpha)Q([s_t, m_t], a_t)$

Given the transition dynamics encoded by the Bayes net in this problem, select whether or not this algorithm needs additional modifications, and why.

☐ Yes. Since the transition dynamics differ from an MDP, we need to account for that in our Q-learning update by weighing the samples according to their likelihood of occurring.

☒ No. The samples are generated from the true dynamics and will, after a sufficient number of updates, represent the empirical probabilities of the true transition probabilities.

Q-learning from lecture does not have an explicit transition function in its update formula, since we can assume that the samples that we collect, after a sufficient number of updates, roughly approximate the transition function $P(s_{t+1} | s_t, a_t)$. For the same reason, combined with the fact that the transition dynamics in the Bayes net of this problem are still Markovian in (s_t, m_t) , we do not need to make any additional modifications to the algorithm above to make it behave like Q-learning on standard MDPs.

- (h) If we were to run value iteration, what would the update formula be for the value $V([s_t, m_t])$, the expected sum of discounted rewards acting optimally from (s_t, m_t) ?

For each letter (A), (B), (C), (D), (E) fill in a single entry for the term corresponding to the correct equation. Recall the definition of $T(s_{t+1})$, $T(m_{t+1})$, and $T(s_{t+1}, m_{t+1})$ from earlier subparts of this problem.

$$V([s_t, m_t]) \leftarrow \text{(A) (B) (C) [(D) + (E) } V([s_{t+1}, m_{t+1}])]$$

- (i) [1 pt] (A) ☒ \max_a ☐ \max_{m_t} ☐ \max_{s_t} ☐ \max_{s_t, m_t} ☐ \sum_a ☐ \sum_{s_t} ☐ \sum_{m_t} ☐ \sum_{s_t, m_t} ☐ 1
- (ii) [1 pt] (B) ☐ \sum_a ☒ $\sum_{m_{t+1}}$ ☐ $\sum_{s_{t+1}}$ ☐ γ ☐ 0 ☐ 1
- (iii) [1 pt] (C) ☒ $\sum_{s_{t+1}} T(s_{t+1}, m_{t+1})$ ☐ $T(m_{t+1})$ ☐ $T(s_{t+1})$ ☐ 0 ☐ 1
- (iv) [1 pt] (D) ☒ $R(m_{t+1})$ ☐ $\sum_{t=1}^N R(m_{t+1})$ ☐ 0 ☐ 1
- (v) [1 pt] (E) ☒ γ ☐ 1

We can simply apply the standard value iteration equation is $V([s_t, m_t]) = \max_a \sum_{s_{t+1}, m_{t+1}} T(s_{t+1}, m_{t+1}) [R(m_{t+1}) + \gamma V([s_{t+1}, m_{t+1}])]$,

which we can interpret as an expectation over next states (s_{t+1}, m_{t+1}) of the sum of discounted rewards.

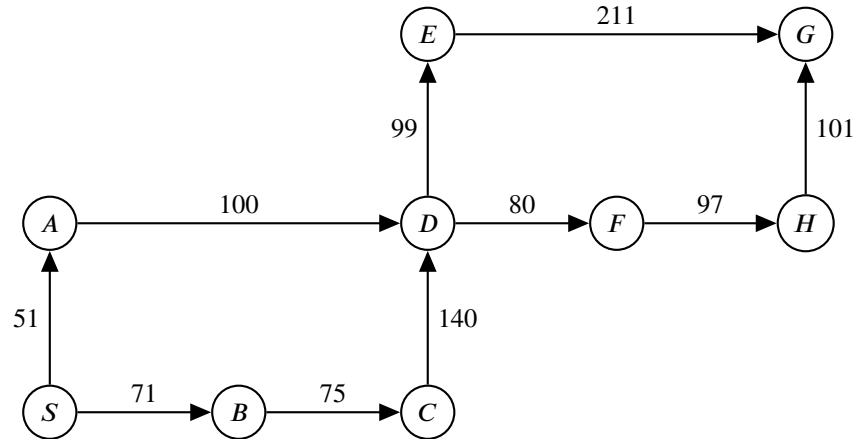
Thus we get the update:

$$\begin{aligned} V([s_t, m_t]) &= \max_a \sum_{s_{t+1}, m_{t+1}} P(s_{t+1}, m_{t+1} | s_t, m_t, a_t) [R(m_{t+1}) + \gamma V([s_{t+1}, m_{t+1}])] \\ &= \max_a \sum_{m_{t+1}} \sum_{s_{t+1}} T(s_{t+1}, m_{t+1}) [R(m_{t+1}) + \gamma V([s_{t+1}, m_{t+1}])] \end{aligned}$$

Q6. [14 pts] Holiday Planning

After a stressful exam period, you plan to spend some time traveling in your new programmable self-driving car. Given two points on the map and a search procedure, the car calculates the optimal route from start to finish.

(a) Let's first frame this problem as a search problem. Consider the following directed graph:



The edges of the graph indicate the cost of the path between the two nodes connected by that edge. We start at node S and we want to reach node G . For the following sub-questions, assume ties resolve in such a way that states with earlier alphabetical order are expanded first.

(i) [1 pt] We first decide to use the graph-search version of Breadth-First Search (BFS) to find the path from S to G . What is the correct order in which states are going to be expanded?

☐ S, A, D, E, G

☒ S, A, B, D, C, E, F, G

☐ S, A, B, D, E, G

☐ S, A, D, F, H, G

When we expand S , the fringe is $\{A, B\}$. We first pop A from the fringe, and append D . Now the fringe is $\{B, D\}$. We then pop B from the fringe and append C . Now the fringe is $\{D, C\}$ so D is popped from the fringe first.

(ii) [2 pts] Now let's take into account the path costs between nodes by running Uniform Cost Search (UCS). What is the *final path* that is going to be returned by UCS?

☐ S, A, D, E, G

☐ S, A, B, D, C, E, F, G

☐ S, A, B, D, E, G

☒ S, A, D, F, H, G

The UCS algorithm finds the shortest path from S to G , which is composed of the shortest path from S to D and the shortest path from D to G , i.e. **SADFHG**.

(iii) [2 pts] Imagine that you are trying to implement informed search for a graph similar to the one above. Which of the following statements are true? Select all that apply.

☒ A heuristic function that is consistent must also be admissible.

☐ The graph-search version of A^* search is optimal, as long as our heuristic function is admissible.

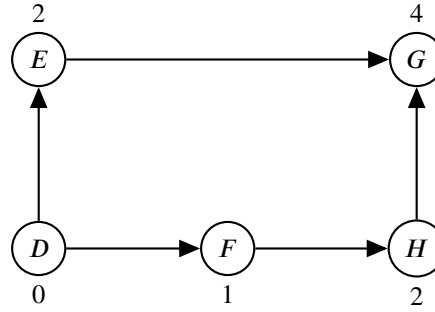
☐ A^* search with a heuristic function $h(n) = 0$ for every node n in the graph expands fewer nodes than UCS.

☒ A non-negative heuristic that never overestimates the cost to reach the goal is admissible.

☐ None of the above

Consistency implies admissibility. For A-star graph search we need consistency. If the heuristic is 0 then A-star is equivalent to UCS. A non-negative heuristic that never overestimates the cost to reach the goal is admissible (this is the definition).

- (b) To account for probabilistic events, we will now frame the problem as a Markov Decision Process (MDP). We will focus on a subset of the nodes from the graph above:



We start at node D and we want to reach node G . From D we can go either to node F or E . We denote the action of moving from D to F as *Move to F* and the action of moving from D to E as *Move to E*. For other nodes, we say that we *Move* when we go from the node to its neighbor in the direction of the edge. Specifically, we can *Move* from F to H , from H to G , and from E to G .

- (i) [2 pts] For actions *Move to F*, *Move to E*, and *Move*, we arrive at the destination node with a probability of 75%, otherwise we transition to a special state called *Broke*. Which of the following values for the transition function are accurate? Select all that apply.

- ☐ $T(D, \text{Move}, s) = \frac{1}{4}$, for $s \in \{E, F\}$
☒ $T(s, \text{Move}, s') = \frac{3}{4}$, for $(s, s') \in \{(F, H), (H, G), (E, G)\}$
☒ $T(s, \text{Move}, \text{Broke}) = \frac{1}{4}$, for $s \in \{D, E, F, H\}$
☒ $T(D, \text{Move}, G) = 0$
☐ None of the above

The transition probability to the *Broke* state is $1/4$, and the transition probability is $3/4$ to the desired state. The third option is considered correct either way, since there is a typo in the problem (the action at D should be "Move to E/F" instead of "move")

- (ii) [2 pts] At any point we can choose to take the action *Stop*, which will transition us to the special state *Done* and yield the reward indicated next to the node in the graph above. For instance, if we take the action *Stop* at node F , we obtain a reward of 1. Which of the following values for the reward function are accurate? Select all that apply.

- ☒ $R(G, \text{Stop}, \text{Done}) = 4$
☒ $R(s, \text{Move}, s') = 0$, for $(s, s') \in \{(F, H), (H, G), (E, G)\}$
☐ $R(D, \text{Move to E}, E) = 2$
☒ $R(s, \text{Stop}, \text{Done}) = 2$, for $s \in \{E, H\}$
☐ None of the above

Only the stop action can incur non-zero reward. Stopping at G incurs reward of 4, and stopping at E or H incurs reward of 2.

- (iii) [3 pts] Now recall the policy improvement equation:

$$\forall s, \pi_{i+1}(s) \leftarrow \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

We begin by performing value iteration given the initial policy π_0 :

States	D	F	H	E	G
π_0	<i>Move to F</i>	<i>Move</i>	<i>Stop</i>	<i>Move</i>	<i>Stop</i>
V^{π_0}	1.125	1.5	2	3	4

What is the updated policy π_1 for each of the following states, given the value of the initial policy above? Use the discount factor $\gamma = 1$.

States	D	F	H	E	G
π_1	<i>Move to E</i>	<i>Move</i>			<i>Stop</i>

For G: the only action is to stop. For F: the value of move is 1.5 while the value of stop is 1. For D: the value of stop is 0, the value of move to E is 2.25, and the value of move to F is 1.125.

(iv) [2 pts] Consider the impact of the discount factor γ on the resulting policy. How would the policy π_1 change if we set the discount factor to $\gamma = 0$?

- ☐ $\pi(D) = \text{Move to F}, \pi(F) = \text{Move}, \pi(H) = \text{Move}, \pi(E) = \text{Move}, \pi(G) = \text{Stop}$
- ☒ $\pi(D) = \text{Stop}, \pi(F) = \text{Stop}, \pi(H) = \text{Stop}, \pi(E) = \text{Stop}, \pi(G) = \text{Stop}$
- ☐ $\pi(D) = \text{Move to F}, \pi(F) = \text{Move}, \pi(H) = \text{Stop}, \pi(E) = \text{Move}, \pi(G) = \text{Stop}$
- ☐ $\pi(D) = \text{Move to E}, \pi(F) = \text{Move}, \pi(H) = \text{Move}, \pi(E) = \text{Move}, \pi(G) = \text{Stop}$

The accumulated discounted reward can only be 0 if we do not stop at the first time step.