

- You have approximately 170 minutes.
- The exam is closed book, closed calculator, and closed notes except your two-page sheet of note.
- Mark your answers ON THE EXAM ITSELF. If you are not sure of your answer you may wish to provide a *brief* explanation. All short answer sections can be successfully answered in a few sentences AT MOST.
- For multiple choice questions,
 - ☐ means mark **all options** that apply
 - ☐ means mark a **single choice**
 - When selecting an answer, please fill in the bubble or square **completely** (● and ■)

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

Your Discussion TA (fill in all that apply):

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

For staff use only:

Q1.	They Hatin', Patrolling (Search Formulation)	/17
Q2.	Search	/17
Q3.	M-O-D-E Pruning	/19
Q4.	LQR	/15
Q5.	Approximate Q-Learning	/12
Q6.	Probability and Bayes Net Representation	/23
Q7.	Decode Your Terror (HMM)	/18
Q8.	Raisins Again (VPI)	/12
Q9.	Machine Learning: Potpourri	/19
Q10.	Perceptrons and Naive Bayes	/24
Q11.	Neural Network	/24
	Total	/200

THIS PAGE IS INTENTIONALLY LEFT BLANK

Q1. [17 pts] They Hatin', Patrolling (Search Formulation)

Recall that in Midterm 1, Pacman bought a car, was speeding in Pac-City, and the police weren't able to catch him. Now Pacman has run out of gas, his car has stopped, and he is currently hiding out at an undisclosed location.

In this problem, you are on the police side, tryin' to catch Pacman!

There are still p police cars in the Pac-city of dimension m by n . In this problem, **all police cars can move, with two distinct integer controls: throttle and steering, but Pacman has to stay stationary**. Once one police car takes an action which lands him in the same grid as Pacman, Pacman will be arrested and the game ends.

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

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

- (a) Suppose you can **only control 1 police car**, and have absolutely no information about the remainder of $p - 1$ police cars, or where Pacman stopped to hide. Also, the police cars can travel up to 6 grid/s so $0 \leq v \leq 6$ at all times.

- (i) [4 pts] What is the **tightest upper bound** on the size of state space, if your goal is to use search to plan a sequence of actions that guarantees Pacman is caught, no matter where Pacman is hiding, or what actions other police cars take. Please note that your state space representation must be able to represent **all** states in the search space.

- (ii) [3 pts] What is the maximum branching factor? Your answer may contain integers, m, n .

- (iii) [2 pts] Which algorithm(s) is/are guaranteed to return a path passing through all grid locations on the grid, if one exists?

- | | |
|---|---|
| <input type="checkbox"/> Depth First Tree Search | <input type="checkbox"/> Breadth First Tree Search |
| <input type="checkbox"/> Depth First Graph Search | <input type="checkbox"/> Breadth First Graph Search |

- (iv) [2 pts] Is Breadth First Graph Search guaranteed to return the path with the shortest number of **time steps**, if one exists?

- ☐ Yes ☐ No

- (b) Now let's suppose you can control **all** p police cars at the same time (and know all their locations), but you still have no information about where Pacman stopped to hide

- (i) [3 pts] Now, you still want to search a sequence of actions such that the paths of p police car combined **pass through all $m * n$ grid locations**. Suppose the size of the state space in part (a) was N_1 , and the size of the state space in this part is N_p . Please select the correct relationship between N_p and N_1

- ☐ $N_p = p * N_1$ ☐ $N_p = p^{N_1}$ ☐ $N_p = (N_1)^p$ ☐ None of the above

- (ii) [3 pts] Suppose the maximum branching factor in part (a) was b_1 , and the maximum branching factor in this part is b_p . Please select the correct relationship between b_p and b_1

- ☐ $b_p = p * b_1$ ☐ $b_p = p^{b_1}$ ☐ $b_p = (b_1)^p$ ☐ None of the above

Q2. [17 pts] Search

For this problem, assume that all of our search algorithms use tree search, unless specified otherwise.

- (a) For each algorithm below, indicate whether the path returned after the modification to the search tree is guaranteed to be identical to the unmodified algorithm. Assume all edge weights are non-negative before modifications.

- (i) [3 pts] Adding additional cost $c > 0$ to every edge weight.

	Yes	No
BFS	<input type="radio"/>	<input type="radio"/>
DFS	<input type="radio"/>	<input type="radio"/>
UCS	<input type="radio"/>	<input type="radio"/>

- (ii) [3 pts] Multiplying a constant $w > 0$ to every edge weight.

	Yes	No
BFS	<input type="radio"/>	<input type="radio"/>
DFS	<input type="radio"/>	<input type="radio"/>
UCS	<input type="radio"/>	<input type="radio"/>

- (b) For part (b), two search algorithms are defined to be **equivalent** if and only if they expand the same states in the same order and return the same path. **Assume all graphs are directed and acyclic.**

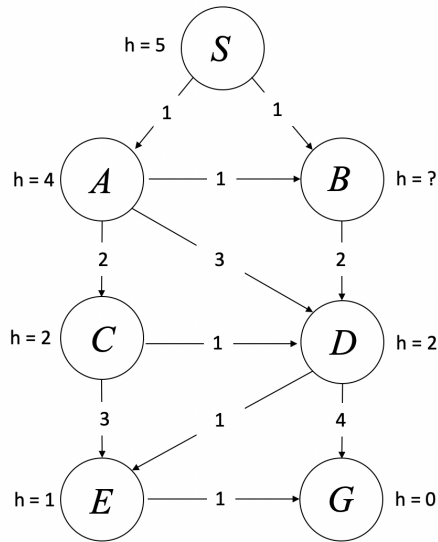
- (i) [3 pts] Assume we have access to costs c_{ij} that make running UCS algorithm with these costs c_{ij} equivalent to running BFS. How can we construct new costs c'_{ij} such that running UCS with these costs is equivalent to running DFS? Mark all correct choices.

- ☐ $c'_{ij} = 0$
☐ $c'_{ij} = 1$
☐ $c'_{ij} = c_{ij}$
☐ $c'_{ij} = -c_{ij}$
☐ $c'_{ij} = c_{ij} + \alpha$
☐ Not possible

- (ii) [3 pts] Given edge weight $c_{ij} = h(j) - h(i)$, where $h(n)$ is the value of the heuristic function at node n , running UCS on this graph is equivalent to running which of the following algorithm(s) on the same graph? Select all correct answers.

- ☐ DFS
 ☐ BFS
 ☐ Iterative Deepening
☐ Greedy
 ☐ A*
 ☐ None of the above.

(c) Consider the following graph. $h(n)$ denotes the heuristic function evaluated at node n .



- (i) [2 pts] Given that G is the goal node, and heuristic values are fixed for all nodes other than B , for which values of $h(B)$ will A* tree search be guaranteed to return the optimal path? Fill in the lower and upper bounds or select “impossible.”

_____ $\leq h(B) \leq$ _____ ☐ Impossible

- (ii) [3 pts] With the heuristic values fixed for all nodes other than B , for which values of $h(B)$ will A* graph search be guaranteed to return the optimal path? Either fill in the lower and upper bound or select “impossible.”

_____ $\leq h(B) \leq$ _____ ☐ Impossible

Q3. [19 pts] M-O-D-E Pruning

For all parts of this question, write your final answer in the small box, but we encourage you to show relevant work in the big box. A correct final answer is sufficient for full credit.

NOTE: In the pruning process, we are only concerned about returning the correct value of the root node.

(a) Minimax Tree Pruning

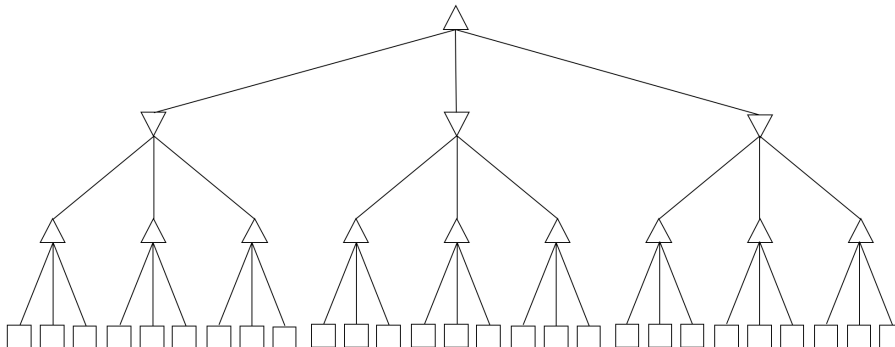


Figure 1: Minimax Tree with $b = 3, d = 3$

Suppose we have a complete minimax tree (with exactly 1 maximizer and 1 minimizer, alternating), where all nodes have branching factor b , and d total layers (an extra layer of maximizer is added, if d is not even)

- (i) [2 pts] In the case of $b = 3, d = 3$, what is the total number of **value nodes** (denoted by squares at the lowest layer)? Your final answer should be an integer.

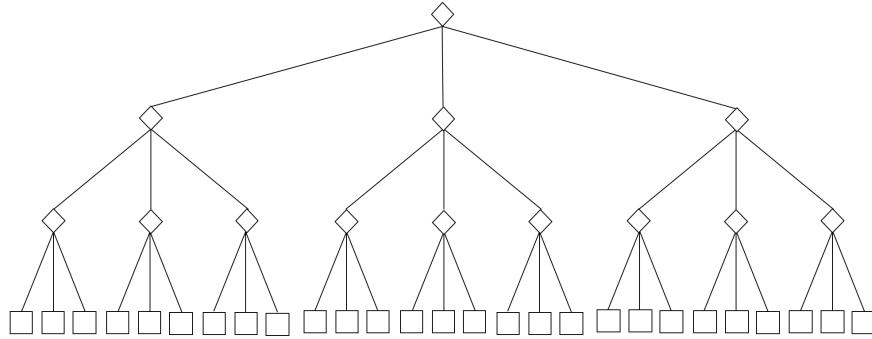
Answer:

- (ii) [3 pts] In the case of $b = 3, d = 3$, what is the **maximum** total number of value nodes **whose value is never explored because an upstream branch is pruned** in one single set of values for the values nodes. Your final answer should be an integer.

Answer:

- (iii) [3 pts] Now we are changing the branching factor b , keeping $d = 3$. In the case of $b = 5, d = 3$, what is the **maximum** total number of value nodes **whose value is never explored because an upstream branch is pruned** in one single set of values for the value nodes. Your final answer should be an integer.

Answer:

(b) Mode Tree PruningFigure 2: Mode Tree with $b = 3, d = 3$

We all know α - β pruning on minimax tree reduces the number of nodes explored to minimum, but what if we change all the maximizer and minimizer nodes to **mode nodes** (symbolized by spades in the diagram), where mode value (most frequent value) among the child nodes will be selected?

Important: Please note that, just like in $\alpha - \beta$ pruning, a child node n_c can be pruned when its parent node n_p are absolutely sure that the value of n_c will have no influence on the final value of the root node.

- (i) [3 pts] In the case of $b = 3, d = 3$, what is the **maximum** total number of value nodes **whose value is never explored because an upstream branch is pruned** in one single set of values for the values nodes. Your final answer should be an integer.

Answer:

- (ii) [4 pts] Now we are changing the branching factor b , keeping $d = 3$. In the case of $b = 5, d = 3$, what is the **maximum** total number of value nodes **whose value is never explored because an upstream branch is pruned** in one single set of values for the values nodes. Your final answer should be an integer.

Answer:

- (iii) [4 pts] Now we are changing the depth d , keeping $b = 3$. In the case of $b = 3, d = 4$, what is the **maximum** total number of value nodes **whose value is never explored because an upstream branch is pruned** in one single set of values for the values nodes. Your final answer should be an integer.

Answer:

Q4. [15 pts] LQR

For all parts of this question, write your final answer in the small box, but we encourage you to show relevant work in the big box. A correct final answer is sufficient for full credit.

Consider the following deterministic MDP with 1-dimensional continuous states and actions and a finite task horizon:

State Space \mathcal{S} : \mathbb{R}

Action Space \mathcal{A} : \mathbb{R}

Reward Function: $R(s, a, s') = -qs^2 - ra^2$ where $r > 0$ and $q \geq 0$

Deterministic Dynamics/Transition Function: $s' = cs + da$ (i.e., the next state s' is a deterministic function of the action a and current state s)

Task Horizon: $T \in \mathbb{N}$

Discount Factor: $\gamma = 1$ (no discount factor)

Hence, we would like to maximize a quadratic reward function that rewards small actions and staying close to the origin. In this problem, we will design an optimal agent π_t^* and also solve for the optimal agent's value function V_t^* for all timesteps.

By induction, we will show that V_t^* is quadratic. Observe that the base case $t = 0$ trivially holds because $V_0^*(s) = 0$. For all parts below, assume that $V_t^*(s) = -p_t s^2$ (Inductive Hypothesis).

- (a) (i) [5 pts] Write the equation for $V_{t+1}^*(s)$ as a function of s, q, r, a, c, d , and p_t . If your expression contains *max*, you do not need to simplify the *max*.

$$V_{t+1}^*(s) =$$

- (ii) [5 pts] Now, solve for $\pi_{t+1}^*(s)$. Recall that you can find local maxima of functions by computing the first derivative and setting it to 0.

$$\pi_{t+1}^*(s) =$$

(b) [5 pts] Assume $\pi_{t+1}^* = k_{t+1}s$ for some $k_{t+1} \in \mathbb{R}$. Solve for p_{t+1} in $V_{t+1}^*(s) = -p_{t+1}s^2$.

$p_{t+1} =$

Q5. [12 pts] Approximate Q-Learning

For all parts of this question, write your final answer in the small box, but we encourage you to show relevant work in the big box. A correct final answer is sufficient for full credit.

We would like to train a Q function Q_w parameterized by some trainable weights w . On observing a sample (s, a, r, s') from the environment, we compare the old estimate of the Q-value of (s, a) : $\hat{y} = Q_w(s, a)$ and the new estimate: $y = r + \max_{a'} Q_w(s', a')$ and use it to update w to minimize the squared error via gradient descent.

$$\begin{aligned} \text{error}(\hat{y}, y) &= (\hat{y} - y)^2 \\ &= (Q_w(s, a) - y)^2 \end{aligned}$$

Treating the target y as a constant, differentiate $\text{error}(\hat{y}, y)$ with respect to w to derive the weight update rule for w_i :

$$w_i \leftarrow w_i - \alpha \frac{\partial \text{error}}{\partial w_i}(\hat{y}, y)$$

Let's derive the update rules for different types of Q functions. For all parts, leave your answers in terms of $y, \hat{y}, f_i(s, a)$, and w_i .

- (a) [4 pts] The Q-function we would like to train has the form $Q_w(s, a) = \sum_{i=0}^m w_i f_i(s, a)$. Please derive $\frac{\partial \text{error}}{\partial w_i}(\hat{y}, y)$ for this Q function.

$$\frac{\partial \text{error}}{\partial w_i}(\hat{y}, y) =$$

(b) The Q-function we would like to train now has the form

$$Q_w(s, a) = \sum_{i=0}^m e^{w_i f_i(s, a)} + \sum_{i=m+1}^{2m} w_i^2 f_i(s, a)$$

(i) [4 pts] For $i \in \{1, \dots, m\}$, derive $\frac{\partial error}{\partial w_i}(\hat{y}, y)$ for this Q function.

$$\frac{\partial error}{\partial w_i}(\hat{y}, y) =$$

(ii) [4 pts] For $i \in \{m+1, \dots, 2m\}$, derive $\frac{\partial error}{\partial w_i}(\hat{y}, y)$ for this Q function.

$$\frac{\partial error}{\partial w_i}(\hat{y}, y) =$$

Q6. [23 pts] Probability and Bayes Net Representation

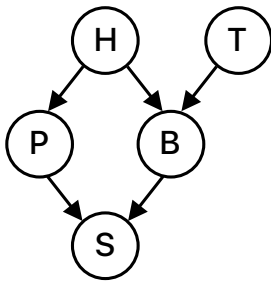
You're interested in knowing whether you would be **S**atisfied with your choice of snack(s), and so you decide to make the prediction using probabilistic inference over a model with the following variables:

- S , whether or not you will be **S**atisfied.
- H , whether or not you will be **H**ungry.
- T , whether or not you will be **T**hirsty.
- P , whether or not you will have **P**izza.
- B , whether or not you will have **B**oba.

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?

- (b) [3 pts] 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.)

- (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.

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

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

$$P(+s|-h,+t) =$$

- (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?

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

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

$$\begin{array}{ccccccc} \bigcirc & f_1(P) & \bigcirc & f_1(B) & \bigcirc & f_1(T) & \bigcirc & f_1(+s) \\ \bigcirc & f_1(P, B) & \bigcirc & f_1(P, T) & \bigcirc & f_1(P, +s) & \bigcirc & f_1(B, T) & \bigcirc & f_1(B, +s) & \bigcirc & f_1(T, +s) \\ \bigcirc & f_1(P, B, T) & \bigcirc & f_1(P, B, +s) & \bigcirc & f_1(B, T, +s) & & & & & & \end{array}$$

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

$$f_1(\text{_____}) = \boxed{\hspace{10cm}}$$

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

--	--

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

$$f_2(\text{_____}) = \boxed{\text{_____}}$$

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

--

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

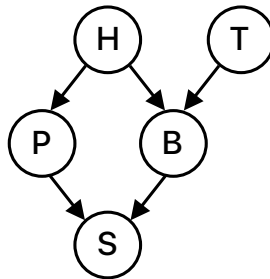
$$f_3(\text{_____}) = \boxed{\hspace{10cm}}$$

- (vi) [2 pts] Write out the expression for computing $\Pr(+b | +s)$ in terms of the remaining factor(s) (after P is eliminated).

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

- (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

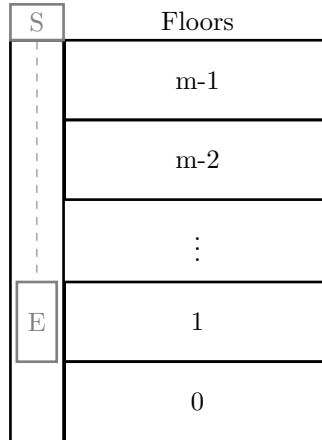
- (vi) [1 pt] $P \perp\!\!\!\perp T | H, S$

☐ Guaranteed

☐ Not guaranteed

Q7. [18 pts] Decode Your Terror (HMM)

You go to Disney and ride the famous Tower of Terror ride, where an elevator rises and drops seemingly at random. You're terrified, but vow to determine the sequence of rises and drops that make up the ride so you won't be as terrified next time. Assume the elevator E follows a Markovian process and it has m floors at which it can stop. In the dead of night, you install a sensor S at the top of the shaft that gives approximate distance measurements, quantized into n different distance bins. Assume that the elevator stops at T floors as part of the ride and the initial distribution of the elevator is uniform over the m floors.



You want to know the most probable sequence of (hidden) states $X_{1:T}$ given your observations $y_{1:T}$ from the sensor, so you turn to the Viterbi algorithm, which performs the following update at each step:

$$m_t[x_t] = P(y_t|x_t) \max_{x_{t-1}} [P(x_t|x_{t-1})m_{t-1}[x_{t-1}]]$$

$$a_t[x_t] = \arg \max_{x_{t-1}} m_{t-1}[x_{t-1}]$$

- (a) (i) [2 pts] What is the run time of the Viterbi algorithm to determine all previous states for this scenario? Please answer in big O notation, in terms of T , m , and n , or write “N/A” if the run time is unable to be determined with the given information.

- (ii) [2 pts] What is the space complexity of the Viterbi algorithm to determine all previous states for this scenario? Please answer in big O notation, in terms of T , m , and n , or write “N/A” if the space complexity is unable to be determined with the given information.

Eventually, we decide that the end of the ride is the exciting part, so we decide that we only wish to determine the previous K states.

- (b) (i) [2 pts] What is the run time of the Viterbi algorithm to determine the previous K states? Please answer in big O notation, in terms of T , K , m , and n , or write “N/A” if the run time is unable to be determined

with the given information.

- (ii) [2 pts] What is the space complexity of the Viterbi algorithm to determine the previous K states? Please answer in big O notation, in terms of T , K , m , and n , or write “N/A” if the space complexity is unable to be determined with the given information.

Suppose you instead only wish to determine the current distribution (at time T) for the elevator, given your T observations, so you use the forward algorithm, with update step shown here:

$$P(X_t|y_t) \propto P(y_t|x_t) \sum_{x_{t-1}} P(X_t|x_{t-1})P(x_{t-1}, y_{0:t-1})$$

Additionally, from your previous analysis, you note that there are some states which are unreachable from others (e.g., the elevator cannot travel from the top floor to the bottom in a single timestep). Specifically, from each state, there are between $G/2$ and G states which can be reached in the next timestep, where $G < m$.

- (c) (i) [2 pts] What is the run time for the forward algorithm to estimate the current state at time T , assuming we ignore states that cannot be reached in each update? Please answer in big O notation, in terms of T , m , G , and n , or write “N/A” if the run time is unable to be determined with the given information.

- (ii) [2 pts] What is the space complexity for the forward algorithm to estimate the current state at time T , assuming we ignore states that cannot be reached in each update? Please answer in big O notation, in terms of T , m , G , and n , or write “N/A” if the space complexity is unable to be determined with the given information.

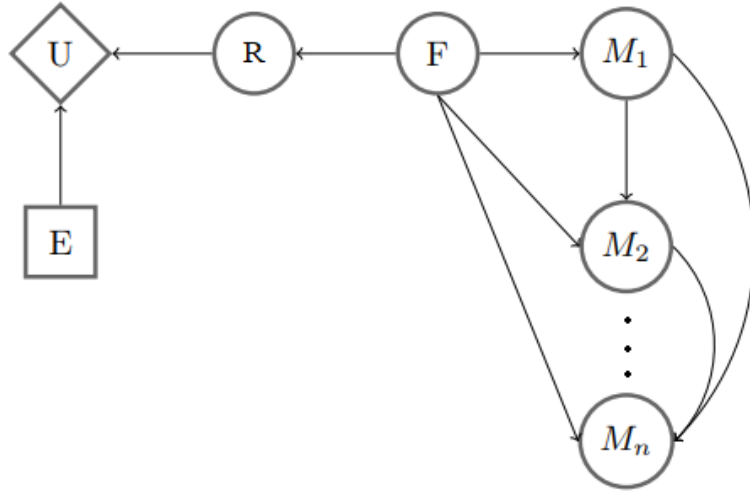
Finally, assume that the number of elevator states is actually infinite (e.g., instead of stopping at floors, the elevator can stop at any point along the elevator shaft).

- (d) [3 pts] What is the run time of the standard forward algorithm to estimate the current state at time T in this case? Please answer in big O notation, in terms of T , m , G , and n , or write “N/A” if the run time is unable to be determined with the given information.

- (e) [3 pts] Suppose you decide to use a particle filter instead of the forward algorithm. What is the run time of a particle filter with P particles? Please answer in big O notation, in terms of T , m , G , n , and P , or write “N/A” if the run time is unable to be determined with the given information.

Q8. [12 pts] Raisins Again (VPI)

- (a) Valerie from midterm 1 is still concerned about her cookie containing raisins so she decides to take a completely new approach. She wants to find out which factory (F) made the cookie because that can help her find out if there are raisins (R) in the cookie. Thus she goes to the cookie company and asks a manager, M_1 , which factory the cookie was made in. However, she doesn't fully trust his answer so she asks his superior, M_2 , which factory he thinks it is in and how credible M_1 is. She knows there is a chain of n managers like this where every manager can tell her which factory they think the cookie was made in and the credibility of **every** manager working under them. This system can be modeled by the decision network below.



For this question assume $1 < i < k < n$, choose one for each equation:

	Could be true	Must be true	Must be false
$VPI(M_i) > VPI(M_k)$	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
$VPI(M_k) > VPI(M_i)$	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
$VPI(F) \geq VPI(M_i)$	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
$VPI(F) > VPI(R)$	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
$VPI(M_k M_i) > VPI(M_k)$	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
$VPI(M_i M_k) > VPI(M_i)$	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
$VPI(M_i, M_k) = VPI(M_i M_k) + VPI(M_k)$	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
$VPI(M_i F) > 0$	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Q9. [19 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.

- (b) [3 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.

- (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

- (j) (i) [2 pts] You train a linear classifier on 1,000 training points and discover that the training accuracy is only 50%. Which of the following, if done in isolation, has a good chance of improving your training accuracy?

☐ Add novel features

☐ Train on more data

☐ Train on less data

- (ii) [2 pts] You now try training a neural network but you find that the training accuracy is still very low. Which of the following, if done in isolation, has a good chance of improving your training accuracy?

☐ Add more hidden layers

☐ Add more units to the hidden layers

- (k) Recall the following kernels covered in class:

1. Linear kernel: $K(x, x') = x \cdot x' = \sum_i x_i x'_i$

2. Quadratic kernel: $K(x, x') = (x \cdot x' + 1)^2 = \sum_{i,j} x_i x_j x'_i x'_j + 2 \sum_i x_i x'_i + 1$

3. Gaussian RBF kernel: $K(x, x') = \exp\left(-\frac{1}{2\sigma^2} \|x - x'\|^2\right)$

- (i) [1 pt] There exists data that is separable with a Gaussian RBF kernel but not with a quadratic kernel.

☐ True

☐ False

- (ii) [1 pt] There exists data that is separable with a linear kernel but not with a quadratic kernel.

☐ True

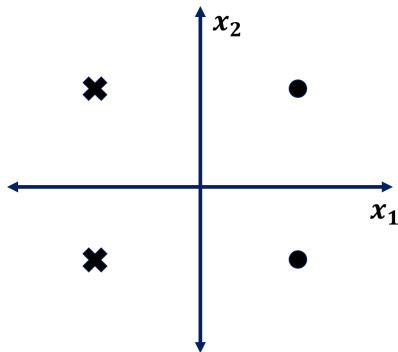
☐ False

Q10. [24 pts] Perceptrons and Naive Bayes

- (a) For each of the datasets represented by the graphs below, please select the feature maps for which the perceptron algorithm can perfectly classify the data.

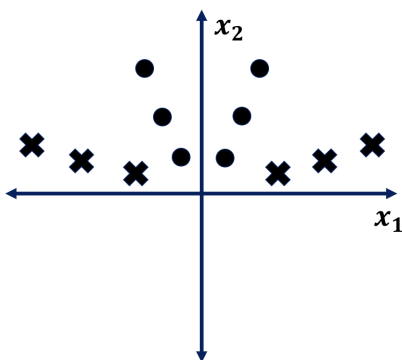
Each data point is in the form (x_1, x_2) , and has some label Y , which is either a 1 (dot) or -1 (cross).

(i) [5 pts]



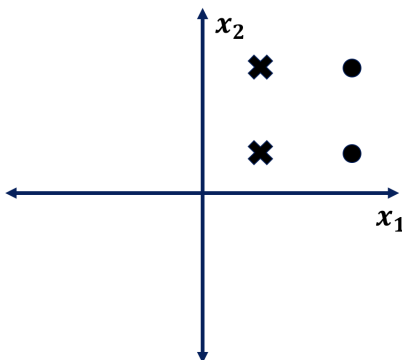
- ☐ $[x_1 \ x_2 \ 1]$
- ☐ $[x_1 \ x_2 \ x_1^2]$
- ☐ $[x_1 \ x_2 \ |x_1|]$
- ☐ $[x_1 \ x_2 \ Y]$
- ☐ $[x_1 \ x_2]$

(ii) [5 pts]



- ☐ $[x_1 \ x_2 \ 1]$
- ☐ $[x_1 \ x_2 \ x_1^2]$
- ☐ $[x_1 \ x_2 \ |x_1|]$
- ☐ $[x_1 \ x_2 \ Y]$
- ☐ $[x_1 \ x_2]$

(iii) [5 pts]



- ☐ $[x_1 \ x_2 \ 1]$
- ☐ $[x_1 \ x_2 \ x_1^2]$
- ☐ $[x_1 \ x_2 \ |x_1|]$
- ☐ $[x_1 \ x_2 \ Y]$
- ☐ $[x_1 \ x_2]$

- (b) [2 pts] Performing maximum likelihood estimation (MLE) to fit the parameters of a Bayes net to some given data (with no Laplace smoothing) leads to which of the following learning algorithms?

- ☐ Naive Bayes
☐ Perceptrons
☐ Kernelization
☐ Neural Networks
☐ None

- (c) Suppose that we are trying to perform a binary classification task using Naive Bayes. Y is the label, and (X_1, X_2) are the features. The domain for the features is anywhere on the 3×3 grid centered at $(0, 0)$. In other words, X_1 and X_2 have the domain $\{-1, 0, 1\}$

Suppose that this is your dataset: $(0, 1, +)$, $(0, -1, -)$, $(-1, 1, +)$, $(-1, -1, -)$, $(1, 0, +)$, $(-1, 1, -)$, $(0, 0, +)$. What is the learned value of each of the following? (Leave your answer as a simplified fraction)

(i) [1 pt]

$$P(Y = +)$$

(ii) [1 pt]

$$P(X_1 = 1 | Y = -)$$

(iii) [1 pt]

$$P(X_2 = 0 | Y = +)$$

- (d) Now, to decouple from the previous question, assume that the learned CPTs are below.

Y	X_1	$\Pr(X_1 Y)$
+	-1	0.4
+	0	0.1
+	1	0.5
-	-1	0.6
-	0	0.3
-	1	0.1

Y	X_2	$\Pr(X_2 Y)$
+	-1	0.2
+	0	0.2
+	1	0.6
-	-1	0.7
-	0	0.1
-	1	0.2

Y	$\Pr(Y)$
+	0.2
-	0.8

- (i) [2 pts] What would be the predicted value for Y if the data point is at $(0, 0)$?

- (ii) [2 pts] What would be the predicted value for Y if the data point is at $(1, -1)$?

Q11. [24 pts] Neural Network

The network below is a neural network with inputs x_1 and x_2 , and outputs y_1 and y_2 . The internal nodes are computed below. All variables are scalar values. Note that $ReLU(x) = \max(0, x)$.

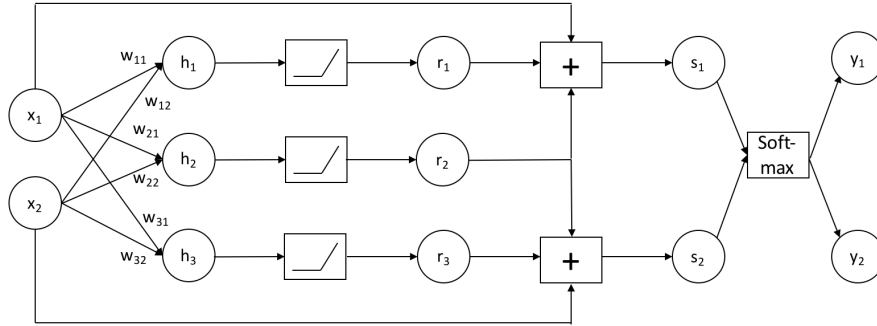


Figure 3: Neural Network

The expressions for the internal nodes in the network are given here for convenience:

$$h_1 = w_{11}x_1 + w_{12}x_2 \quad h_2 = w_{21}x_1 + w_{22}x_2 \quad h_3 = w_{31}x_1 + w_{32}x_2$$

$$r_1 = ReLU(h_1) \quad r_2 = ReLU(h_2) \quad r_3 = ReLU(h_3) \quad s_1 = x_1 + r_1 + r_2 \quad s_2 = x_2 + r_2 + r_3$$

$$y_1 = \frac{\exp(s_1)}{\exp(s_1) + \exp(s_2)} \quad y_2 = \frac{\exp(s_2)}{\exp(s_1) + \exp(s_2)}$$

(a) Forward Propagation

Suppose for this part only, $x_1 = 3, x_2 = 5, w_{11} = -10, w_{12} = 7, w_{21} = 2, w_{22} = 5, w_{31} = 4, w_{32} = -4$. What are the values of the following internal nodes? Please simplify any fractions.

(i) [1 pt] $h_1 =$

(ii) [2 pts] $s_1 =$

(iii) [3 pts] $y_2 =$

(b) Back Propagation

Compute the following gradients analytically. The answer should be an expression of any of the nodes in the network ($x_1, x_2, h_1, h_2, h_3, r_1, r_2, r_3, s_1, s_2, y_1, y_2$) or weights $w_{11}, w_{12}, w_{21}, w_{22}, w_{31}, w_{32}$ (clarification during exam: without derivative or partial derivative symbols). In the case where the gradient depend on the value of nodes in the network, **please list all possible analytical expressions, caused by active/inactive ReLU, separated by comma.**

Hint 1: If z is a function of y , and y is a function of x , the chain rule of taking derivative is: $\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} * \frac{\partial y}{\partial x}$

Hint 2: Hint: Recall that for functions of the form $g(x) = \frac{1}{1 + \exp(a - x)}$, $\frac{\partial g}{\partial x} = g(x)(1 - g(x))$

(i) [1 pt] $\frac{\partial h_2}{\partial x_1} =$

(ii) [2 pts] $\frac{\partial h_1}{\partial w_{21}} =$

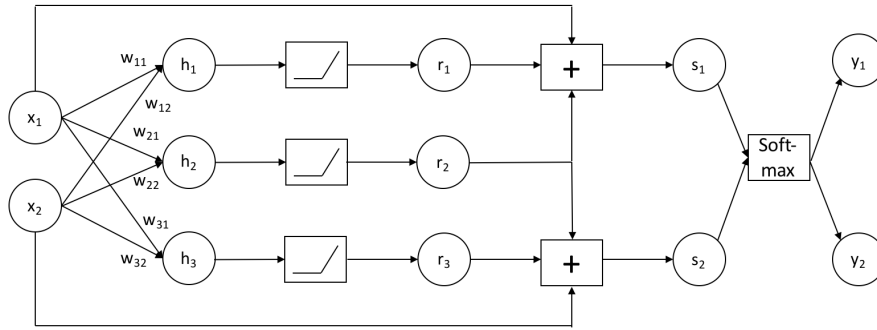


Figure 4: Neural Network, copied from the previous page for reference

(iii) [2 pts] $\frac{\partial r_3}{\partial w_{31}} =$

(iv) [2 pts] $\frac{\partial s_1}{\partial r_1} =$

(v) [3 pts] $\frac{\partial s_1}{\partial x_1} =$

(vi) [3 pts] $\frac{\partial y_2}{\partial s_2} =$

(vii) [3 pts] $\frac{\partial y_1}{\partial x_1} =$

- (c) [2 pts] In roughly 15 words, what role could the non-negative values in node r_2 play according to its location in the network architecture?