

CSCI 3202
Final Exam
Fall 2020

Write clearly and in the box:

Name: Jacob (Jake) Henson
Student ID: 105963531
Section number: 001

- **RIGHT NOW!** Include your name, student ID and section number on the top of your exam. If you're handwriting your exam, write this information at the top of the first page!
- You may use the textbook, your notes, lecture materials, and Piazza as recourses. Piazza posts should not be about exact exam questions, but you may ask for technical clarifications and ask for help on review/past exam questions that might help you. You may not use external sources from the internet or collaborate with your peers.
- You may use a calculator.
- If you print a copy of the exam, clearly mark answers to multiple choice questions in the provided answer box. If you type or hand-write your exam answers, write each problem on their own line, clearly indicating both the problem number and answer letter. Start each new problem on a new page.
- Mark only one answer for multiple choice questions. If you think two answers are correct, mark the answer that **best** answers the question. No justification is required for multiple choice questions. For handwriting multiple choice answers, clearly mark both the number of the problem and your answer for each and every problem.
- For free response questions you must clearly justify all conclusions to receive full credit. A correct answer with no supporting work will receive no credit.
- The Exam is due to Gradescope by midnight on Monday, December 10.
- When submitting your exam to Gradescope, use their submission tool to mark on which pages you answered specific questions.

Problem	Max Points
Short Response	20
Markov Models	20
Bayes' Nets	20
Learning	20
MDP	20
Total	100

(1) [20 points] Short responses. Provide justification when asked.

1A) [5 points] Which of the following five algorithms are guaranteed to find an *optimal* solution to a search problem for the shortest path cost on a given finite graph? **Circle** all that apply.

Breadth-First Search;

Depth-First Search

Greedy Best-First;

A* with *any* heuristic;

Uniform-Cost Search

1B) [5 points] What is the major difference between active and passive learning? Provide an example of each and explain why those examples highlight the difference between the types of learning.

Passive Learning is where the agent learns the utilities of states from a fixed pre-determined policy, whereas Active Learning the agent must learn to update actions/policies along the way. An example of passive learning is with Direct Estimation, where we simply try a whole bunch of policies and pick the best via brute-force, whereas with Active Learning, an example is the multi-armed bandit where we use a greedy strategy to keep estimating best policy, keeping track of this with $Q[k]$.

1C) [5 points] In our course discussion of the ϵ -greedy algorithm, we take the original or best action with probability $1 - \epsilon$ and take a random action drawn from a uniform distribution with probability ϵ . If we can tune or adjust ϵ over the course of a search problem, would we initialize it with a low or high value at the start of training? What about at the end of the training epochs? Answer both prompts and justify your choices.

- At the start of training, I think it is best to set ϵ high so that the agent has a higher probability to explore the environment rather than exploit it early on.

- At the end of training, it is best to set ϵ lower, so that the agent will be more likely exploit the environment and be less likely to explore.

1D) [5 points] True or False, and Justify. In Q-learning all samples must be from the optimal policy to find optimal (correct) q -values.

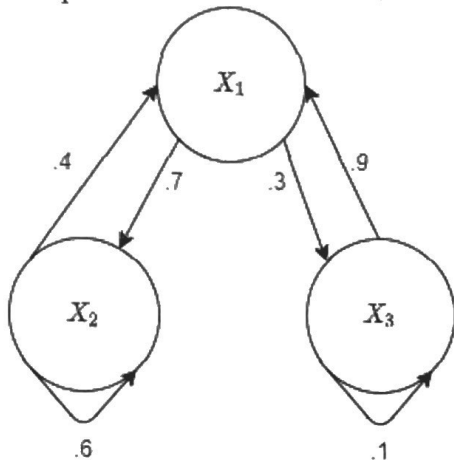
False - a q -value (action, utility) may be optimal from a non-optimal policy if it is the value which is deemed correct as it continues to learn. Eventually, it will converge to an optimal value over time.

(i.e., a q value can give a large return even in a sub-optimal policy)

2) [20 points] Markov Models.

3

- 2A) [7 points] Consider the following graph which represents a Markov model, with state transition probabilities for X shown along each edge. **Set up** (you don't have to solve) a system of equations whose solution is the long-run distribution or *stationary distribution* for X :



$$\begin{aligned}
 P(X_1) &= P(X_1|X_2)P(X_2) + P(X_1|X_3)P(X_3) \\
 P(X_2) &= P(X_2|X_1)P(X_1) \\
 P(X_3) &= P(X_3|X_1)P(X_1) \\
 P(X) &= P(X_1) + P(X_2) + P(X_3)?
 \end{aligned}$$

Parts 2B and 2C refer to the following: With the same 3 states for X as above, suppose you're designing a smart house lighting AI that you've cleverly named Housing Automated Lighting (HAL). HAL's sole purpose is to turn the lights on in whichever room you're in. Your house has the same 3 rooms/states X as above, and the graph above represents your probability of moving from one room to the next each hour of the day. At the *start* of the day, ($t = 0$) you always begin in your bedroom marked X_1 above, which we'll represent as the event $X_0 = 1$.

So that HAL may have enough information to track which room you're in, it's equipped with a sensor Y , which is a little noisy. When you're actually in room i , it senses that $Y = i$ with probability 80%, while it incorrectly diagnoses that you're in each one of the other rooms 10% of the time. In other words, at time j , the sensed room will be $P(Y_j = i | X_j = i) = 0.8$ for the room you're in and $P(Y_j = i | X_j \neq i) = 0.1$ for each other room.

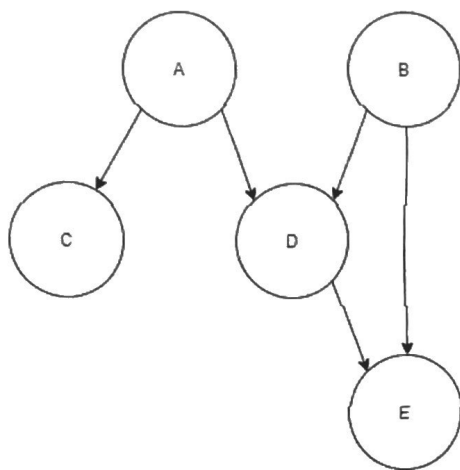
- 2B) [7 points] You wake up as usual at time $t = 0$ and stumble out of bed to begin your day. It's now time $t = 1$. HAL *senses* you in the kitchen (room 2), or $Y_1 = 2$. What is the probability you're actually in the kitchen at time $t = 1$? In other words, what is $P(X_1 = 2)$ given where you started and Hal's sensor?

$$\begin{aligned}
 &P(X_1 = 2 | t=1 \text{ and } Y_1 = 2) \\
 &0.7 \text{ prob of going to kitchen } (X_2) \\
 &0.8 \text{ prob of HAL sensing in kitchen } (Y_2) \\
 &0.56 + (0.3 \cdot 0.1) \text{ or } (P(X_3) \cdot P(Y_3)) = 0.56 + 0.03 = 0.59
 \end{aligned}$$

- 2C) [6 points] After sensing that $Y_1 = 2$, HAL then *again* senses you in the kitchen, so $Y_2 = 2$. Should this increase or decrease the conditional probability that $P(X_1 = 2)$? You may attempt to answer with a fully-explained intuitive answer or perform all exact calculations if you wish.

this conditional probability would probably decrease - you have a 60% chance to stay in the kitchen, whereas originally there was a 70% chance to go, and given that HAL has an 80% chance of getting his sensor guess right, the probability you've stayed in the kitchen is slightly less than the probability you went there in the first place. (i.e. you have less of a chance of BEING IN the kitchen)

- 3) [20 points] Consider the following Bayesian Network, where all variable nodes may only take on True/False as values.



$$P(A = \text{True}) = 0.3$$

$$P(B = \text{True}) = 0.6$$

CPT for C:

A	$P(C = \text{True})$
True	0.6
False	0.3

CPT for D:

A	B	$P(D = \text{True})$
True	True	0.4
True	False	0.5
False	True	0.1
False	False	0.2

CPT for E:

B	D	$P(E = \text{True})$
True	True	0.2
True	False	0.3
False	True	0.5
False	False	0.9

Show all work for the following queries:

- 3A) [7 points] What is the probability that all five variables are simultaneously true?

$$P(A, B, C, D, E) = P(A) \cdot P(B) \cdot P(C|A) \cdot P(E|B, D) \cdot P(D|A, B)$$

$$= 0.3 \cdot 0.6 \cdot 0.6 \cdot 0.4 \cdot 0.2$$

$$= 0.00864$$

- 3B) [7 points] What is the probability that A is false given that all four other variables are true?

$$P(\bar{A} | B, C, D, E) = P(\bar{A}) \cdot P(B) \cdot P(C|\bar{A}) \cdot P(E|B, D) \cdot P(D|\bar{A}, B)$$

$$= 0.7 \cdot 0.6 \cdot 0.3 \cdot 0.2 \cdot 0.1$$

$$= 0.00252$$

- 3C) [6 points] What is the probability that C is true given that D is true?

$$P(C|D) = \frac{P(C, D)}{P(D)} = \frac{0.3 \cdot 0.6 \cdot 0.4 \cdot 0.5 \cdot 0.1 \cdot 0.2}{0.4 \cdot 0.5 \cdot 0.1 \cdot 0.2}$$

$$= \frac{0.00072}{0.0004}$$

$$= 0.18$$

product of events where C and D are true

$$Q_{i+1}(s,a) = Q_i(s,a) + \alpha [R(s) + \max_{a'} Q_i(s',a') - Q_i(s,a)]$$

- 4) [20 points] Suppose an agent exists on state space with three states, X , Y , and Z , which each hold some features about the kittens and puppies currently playing with our agent. Within each state, the agent has two actions: pet the dog (denoted "dog") and pet the cat ("cat"). The agent chooses actions according to policy π , but does not know the underlying process that dictates state transitions. So it sets up an experiment, wherein it:

- Starts at a *random* state s , and chooses an action a .
- Observes the successor state s' of that action and the rewards r resulting from that transition.
- Updates Q -values for that state-action pair.

Suppose that the Q -values are all initialized to 0, the learning rate is fixed as $\alpha = \frac{1}{3}$, and there is no discount factor ($\gamma = 1$). The first 6 training episodes are shown at left below.

- 4A) [12 points] Run Q -learning updates on the table at left, updating the desired quantities to the right in the order of the training episodes:

s	a	s'	r
X	"cat"	Y	3
Z	"dog"	Y	3
Y	"cat"	Z	-3
X	"cat"	Y	6
Y	"dog"	X	0
X	"cat"	Z	3

$$Q_0(s,a) = \alpha [R(s) + \max_{a'} Q_i(s',a') - Q_i(s,a)]$$

$$Q_1(X, \text{cat}) = 0 + \frac{1}{3} [3 + 1 \cdot 0 - 0] = 1.333$$

$$Q_1(Z, \text{dog}) = 0 + \frac{1}{3} [3 + 1 \cdot 0 - 0] = 1.333$$

$$Q_1(Y, \text{cat}) = 0 + \frac{1}{3} [-3 + 1 \cdot 3 - 0] = -4$$

$$Q_2(X, \text{cat}) = 1.333 + \frac{1}{3} [6 + 1 \cdot 3 - 1.333] = 3.888$$

$$Q_1(Y, \text{dog}) = 0 + \frac{1}{3} [0 + 1 \cdot 6 - 0] = 2$$

$$Q_3(X, \text{cat}) = 3.8888 + \frac{1}{3} [3 + 1 \cdot 3 - 3.888] = 4.593$$

- 4B) [4 points] After your training episodes, our agent constructs a policy π that maximizes the estimated utility in a given state. What are the actions chosen by the agent in states X and Y ?

$\pi(s) = (X, \text{cat})$ where possible
since, of course, these values
maximize utility in all examples
above.

- 4C) [4 points] Suppose our agent decides to *estimate* the underlying state transitions from its empirical results in 4A, as it would in *adaptive dynamic* reinforcement learning. Without using any augmented counting such as Laplace smoothing, what would you estimate from the training episodes for:

$$P(s' = Y | s = X, a = \text{"cat"}) = \frac{2}{3}$$

because 2 actions result
in Y from (X, cat)

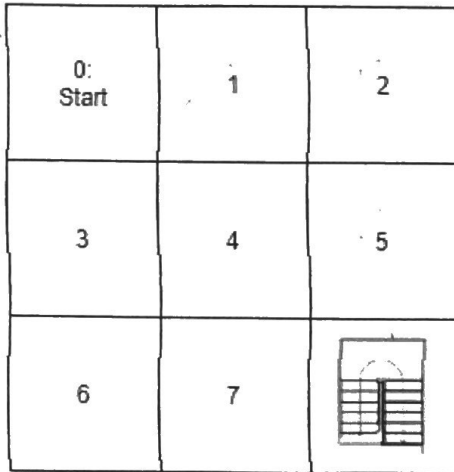
$$P(s' = Z | s = X, a = \text{"cat"}) = \frac{1}{2}$$

because 1 action results
in Z from (X, cat)

$$U_{i+1}(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U_i(s')$$

- 5) [20 points] MDP. Consider the MDP at left below, where an agent starts in a dangerous dungeon. The agent has the standard actions of movement NSEW to non-wall locations. There is a reward of 1 for escaping the floor at the staircase, which represents a terminal state (tile that would be indexed 8). Suppose that the discount factor is $\gamma = 1$ (so no discounting) and there is no reward associated with any state other than the stairs.

Was told
on Piazza
there's no
uncertainty



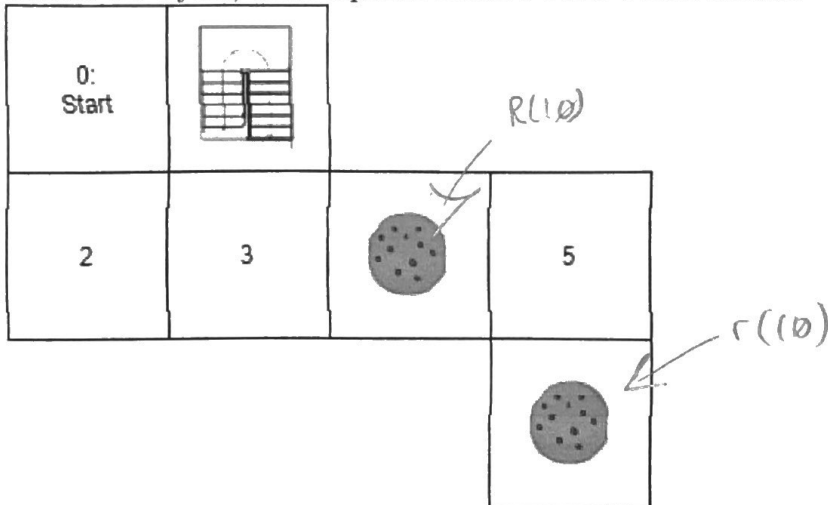
- 5a) [8 points] Complete the following table where each row represents a step of value iteration:

k	U(0)	U(1)	U(2)	U(3)	U(4)	U(5)	U(6)	U(7)	U(8)
0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	1
2	0	0	0	0	0	1	0	1	1
3	0	0	0	0	1	1	0	1	1
4	0	0	1	0	1	1	1	1	1
5	0	1	1	1	1	1	1	1	1
6	1	1	1	1	1	1	1	1	1
7	1	1	1	1	1	1	1	1	1

- 5b) [4 points] Value iteration should have converged in the steps above, such that $U_k(s) = U_{k+1}(s)$ for all states. At which step did it do so?

Step 6. This makes sense because it's a
3x3 matrix,

- 5c) [8 points] Suppose we reach the next level, and now the dungeon floors also are shockingly housing delicious cookies that provide a *one-time* reward of 10 each! You enter the level and observe the layout, where squared indexed 4 and 6 hold cookies:



List *all* elements of the state space and the optimal policies for each assuming that each state now has a (punishment) reward of -0.05 for remaining in the level. The discount factor is still 1. You may assume that the agent can never occupy a space with an uneaten cookie, but you should include states that could only be reached via the agent jumping over or teleporting past a cookie. The first state-optimal policy pair is provided below.

Location(agent)	State(cookie #4)	State(cookie #6)	Action
0	eaten	eaten	East
0	eaten	uneaten	SOUTH
0	uneaten	uneaten	SOUTH
2	eaten	eaten	NORTH
2	eaten	uneaten	EAST
2	uneaten	uneaten	EAST
3	eaten	eaten	NORTH
3	eaten	uneaten	EAST
3	uneaten	uneaten	EAST
4	eaten	eaten	EAST
4	eaten	uneaten	WEST
5	eaten	eaten	WEST
5	eaten	uneaten	SOUTH
5	uneaten	uneaten	SOUTH
6	eaten	eaten	North
6	uneaten	eaten	North

$$r = -0.05$$

$$\gamma = 1$$

You do not need to justify how you arrived at each policy. **Have a great break!**