

Due: Wednesday 03/11/2020 at 11:59pm (submit via Gradescope).

Policy: Can be solved in groups (acknowledge collaborators) but must be written up individually

Submission: Your submission should be a PDF that matches this template. Each page of the PDF should align with the corresponding page of the template (page 1 has name/collaborators, question 1 begins on page 2, etc.). **Do not reorder, split, combine, or add extra pages.** The intention is that you print out the template, write on the page in pen/pencil, and then scan or take pictures of the pages to make your submission. You may also fill out this template digitally (e.g. using a tablet.)

First name	
Last name	
SID	
Collaborators	

For staff use only:

Q1.	MDP: Eating Chocolate	/40
Q2.	MDPs and RL	/30
Q3.	Probability: Flowers	/30
Total		/100

Q1. [40 pts] MDP: Eating Chocolate

We have a chocolate bar of dimensions 1×8 , which contains 8 squares. Most of these squares are delicious chocolate squares, but some of them are poison! Although the chocolate and poison squares are visually indistinguishable, someone has told us which ones are which. The layout for our chocolate bar is shown below with P indicating poison squares.

P				P		P	
---	--	--	--	---	--	---	--

Eating a chocolate square immediately gives a reward of 1 and eating a poison square immediately gives a reward of -2 . **Starting from the right end of the bar**, there are 3 possible actions that we can take (at each step), and all of these actions cause non-deterministic transitions as follows:

- Action b_1 : Try to bite 1 square, which will result in you actually eating 0, 1, or 2 squares with equal probability.
- Action b_2 : Try to bite 2 squares, which will result in you actually eating 1, 2, or 3 squares with equal probability.
- Action *Stop*: Stop biting, which will end the game definitively and result in no reward.

(a) [10 pts] Formulate this problem as an MDP.

States: Decide (and explain) how you want to represent a “state.” Using that representation, list out all of your possible states in this MDP.

Actions:

Transitions:

Rewards: Skip the explicit formulation of the rewards for this problem.

(b) Value Iteration

- (i) [2 pts] Perform value iteration for 3 iterations, using $\gamma = 1$. What are the values for each state, after each of these iterations?

State	1	2	3	4	5	6	7	8
$V_0(s)$	0	0	0	0	0	0	0	0
$V_1(s)$								
$V_2(s)$								
$V_3(s)$								

- (ii) [1 pt] We consider value iteration to have converged by iteration k if $\forall s, V_{k+1}(s) = V_k(s)$. Did the values converge by iteration 2?

☐ Yes ☐ No

- (iii) [1 pt] Given that we know a $V(s)$ for all states, how do we extract a policy $\pi(s)$ from that value function? Your answer should be a one-line math expression which uses some or all of the following: $s, a, s', V, R, T, \gamma, \sum, \text{argmax}, \text{max}$.

$\pi(s) =$

- (iv) [1 pt] What is the resulting policy after extracting it from the values $V_2(s)$? If applicable, write all of the valid actions for a given state.

State	1	2	3	4	5	6	7	8
$\pi_3(s)$								

- (v) [5 pts] Conceptual: In general, we consider a policy to have converged if, for all states s , $\pi(s)$ does not change with further calculation steps. Your friend Victor claims that **in any MDP setting, not limited to this game**, $\pi(s)$ can only be considered to be converged if the $V(s)$ that it came from was already converged. Is Victor right?

- ☐ Victor is right, and I will provide a proof below.
☐ Victor is wrong, and I will provide a counter example below.

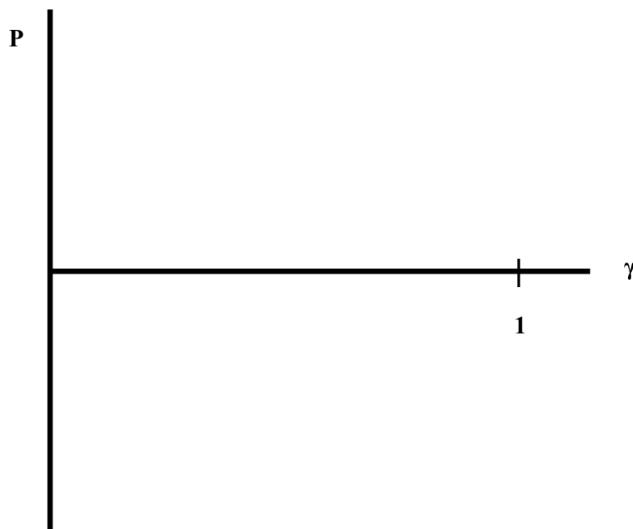
Proof/counter example:

For the rest of the problem, assume we have deterministic transitions. This means that taking the b_1 action bites exactly 1 square and taking the b_2 actions bites exactly 2 squares. Note that if you take b_1 , then your reward for that step is the reward associated with that one square, and if you take b_2 , then your reward for that step is the sum of rewards from the two squares.

- (c) (i) [6 pts] Let's say that you took a sequence of 3 single bites, so you are now in state $s = 5$. At this point, you can either stop or continue to bite. For what range of discount values γ would you choose to stop versus continue to bite? Show your work and explain your answer.

- (ii) [4 pts] For this part, assume the same situation as the previous part, where you took 3 single bites already and you are now in state $s = 5$. Think about how the behavior of the optimal policy at this point changes, as we vary (a) the reward associated with eating a poison square and (b) the discount γ from 0 to 1. First, what's the mathematical condition that must hold for us to choose to continue to bite:

Plot this condition in the plot below by shading in the areas where we will continue to bite. Label important points, including the point where $P = -2$ with the threshold γ value you found in part (c.i).



- (d) (i) [2 pts] Define/create a policy below such that performing value iteration on that policy would give these values indicated here.

State	1	2	3	4	5	6	7	8
$V^\pi(s)$	0	-1	2	3	0	2	-1	0
$\pi(s)$								

- (ii) [3 pts] Is the previous policy optimal? Perform one step of policy iteration to justify your answer. Please show your work.

State	1	2	3	4	5	6	7	8
$\pi_{i+1}(s)$								

Was the policy in part (d.i) optimal: ☐ Yes ☐ No

(iii) [5 pts]

We can build a graph from an MDP as follows. The vertices are the states, and there is an edge from state s to state s' labelled with an action a if $T(s, a, s')$ is nonzero. We can fully represent the MDP and assign a weight of $T(s, a, s')$ and a value of $R(s, a, s')$ to the corresponding edge from s to s' in the graph. Note that there may be more than one edge from s to s' if there is more than one action with a positive transition probability between s and s' . That is, there may be parallel edges in the graph.

Indeed, the MDP about the going fast and slow in a car in lecture is illustrated as a graph in the slides.

Notice that when viewed as a graph the state transitions of the above MDP forms a directed acyclic graph (DAG). It is useful here to explicitly add a terminal state named 0. The terminal state is reached from any state using the S action and every other action from any state leads to a lower numbered state or the terminal state. (To be sure, this would be true as well for the probabilistic version if action b_1 had zero probability of eating nothing.)

Sketch an algorithm to produce the optimal policy for $\gamma = 1$ in time linear in the number of states and edges in the associated graph in an MDP where the transitions form a DAG. You may assume there is a single terminal state with no outgoing transitions. (Hint: the value of the terminal state is 0.) Can you also produce the optimal policy for an arbitrary $\gamma < 1$?

Q2. [30 pts] MDPs and RL

The agent is in a 2×4 gridworld as shown in the figure. We start from square 1 and finish in square 8. When square 8 is reached, we receive a reward of +10 at the game end. For anything else, we receive a constant reward of -1 (you can think of this as a time penalty).

1	2	3	4
5	6	7	8

The actions in this MDP include: up, down, left and right. The agent cannot take actions that take them off the board. In the table below, we provide initial non-zero estimates of Q values (Q values for invalid actions are left as blanks):

Table 1

	action=up	action=down	action=left	action=right
state=1		Q(1, down)=4		Q(1, right)=3
state=2		Q(2, down)=6	Q(2, left)=4	Q(2, right)=5
state=3		Q(3, down)=8	Q(3, left)=5	Q(3, right)=7
state=4		Q(4, down)=9	Q(3, left)=6	
state=5	Q(5, up)=5			Q(5, right)=6
state=6	Q(6, up)=4		Q(6, left)=5	Q(6, right)=7
state=7	Q(7, up)=6		Q(7, left)=6	Q(7, right)=8

- (a) Your friend Adam guesses that the actions in this MDP are fully deterministic (e.g. taking down from 2 will land you in 6 with probability 1 and everywhere else with probability 0). Since we have full knowledge of T and R , we can thus use the Bellman equation to improve (i.e., further update) the initial Q estimates.

Adam tells you to use the following update rule for Q values, where he assumes that your policy is greedy and thus does $\max_a Q(s, a)$. The update rule he prescribes is as follows:

$$Q_{k+1}(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

- (i) [1 pt] Perform one update of $Q(3, \text{left})$ using the equation above, where $\gamma = 0.9$. You may break ties in any way.

- (ii) [3 pts] For the Q update rule prescribed above, how does $Q_{k+1}(s, a)$ depend on $Q_k(s, a)$? How is this different from the normal Q learning update that we saw in lecture, which is $Q_{k+1}(s, a) = (1 - \alpha)Q_k(s, a) + \alpha * \text{sample}$?

- (b) After observing the agent for a while, Adam realized that his assumption of T being deterministic is wrong in one specific way: **when the agent tries to legally move down, it occasionally ends up moving left instead** (except from grid 1 where moving left results in out-of-bound). All other movements are still deterministic.

Suppose we have run the Q updates outlined in the equation above until convergence, to get $Q_{wrong}^*(s, a)$ under the original assumption of the wrong (deterministic) T . Suppose $Q_{correct}^*(s, a)$ denotes the Q values under the new correct T . Note that you don't explicitly know the exact probabilities associated with this new T , but you know that it qualitatively differs in the way described above. As prompted below, list the set of (s, a) pairs where $Q_{wrong}^*(s, a)$ is either an over-estimate or under-estimate of $Q_{correct}^*(s, a)$.

- (i) [3 pts] List of (s, a) where $Q_{wrong}^*(s, a)$ is an over-estimate (and why):

- (ii) [3 pts] List of (s, a) where $Q_{wrong}^*(s, a)$ is an under-estimate (and why):

- (c) [2 pts] Suppose that we have a mysterious oracle that can give us either all the correct Q-values $Q(s, a)$ or all the correct state values $V(s)$. Which one do you prefer to be given if you want to use it to find a greedy policy, and why?
- (d) [2 pts] Suppose that you perform actions in this grid and observe the following episode: 3, right, 4, down, 8 (terminal).
 With learning rate $\alpha = 0.2$, discount $\gamma = 0.9$, perform an update of $Q(3, \text{right})$ and $Q(4, \text{down})$. Note that here, we update Q values based on the sampled actions as in TD learning, rather than the greedy actions.
- (e) [2 pts] One way to encourage an agent to perform more exploration in the world is known as the “ ϵ -greedy” algorithm. For any given policy $\pi(s)$, this algorithm says to take the original action $a = \pi(s)$ with probability $(1 - \epsilon)$, and to take a random action (drawn from a uniform distribution over all legal actions) with probability ϵ . If ϵ can be tuned, would you assign it to be a high or low value at the beginning of training? What about at the end of the training? Please answer both questions and justify your choices.

- (f) Instead of using the “ ϵ -greedy” algorithm, we will now do some interesting exploration with softmax. We first introduce a new type of policy: A stochastic policy $\pi(a|s)$ represents the probability of action a being prescribed, conditioned on the current state. In other words, the policy is now a distribution over possible actions, rather than a function that outputs a deterministic action.

Let's define a new policy as follows:

$$\pi(a|s) = \frac{e^{Q(s,a)}}{\sum_{a'} e^{Q(s,a')}}.$$

- (i) [1 pt] Suppose we are at square 3 in the grid and we want to use the originally provided Q values from the table. What is the probability that this policy will tell us to go right? Note that the sum over actions prescribed above refers to a sum over legal actions.
- (ii) [3 pts] What is the advantage of this exploration strategy, compared with “ ϵ -greedy”? How are they qualitatively different?

- (g) [10 pts] Your friend Cody argues that we could still explicitly calculate Q updates (like Adam's approach in part (a)) even if we don't know the true underlying transition function $T(s, a, s')$, because he believes that our T can be roughly approximated from samples.
- (i) [5 pts] Consider that you are given a moderate amount of (maybe 100,000) of transitions, in the form of (s_{start}, a, s_{end}) . In a few sentences, describe **a sequence of steps** that you can take to compute $T_{approx}(s, a, s')$, which is an approximation of the true underlying (unknown) $T(s, a, s')$.
- (ii) [5 pts] Now, consider a case where Cody does the following: He (1) uses samples to compute $T_{approx}(s, a, s')$ in the way that you described in the previous part, and then (2) uses that $T_{approx}(s, a, s')$ to perform Q updates with Adam's equations. What could be one potential problem with this approach?
Hint: If your $T_{approx}(s, a, s')$ is not exactly correct (i.e., you didn't get infinite data samples, so it's probably not exactly correct), what could go wrong?

Q3. [30 pts] Probability: Flowers

Your friend Ethan wants to find the biggest flower in the field. In the field, there are n flowers where the value of n is unknown. The size of all the flowers can be ranked unambiguously if all are seen. Ethan will walk across this field, with each order of flowers being equally likely. Ethan decides not to walk one meter back because he is too tired. He also only gives him one chance to pick the flower. That means whenever he picks a flower, he is done with the whole sequence of flowers. He wants to have the highest probability of selecting the actual biggest flower.

- (a) (i) [3 pts] Ethan thinks he can do the following: He will check out $r - 1$ flowers without touching them so that he can have a good calibration. M is the biggest flower in the $r - 1$ flowers. He then selects the first subsequent flowers that is better than M . For an arbitrary cutoff r , can you show the **probability that the biggest flower is actually selected**.

The probability

$$P(r) = \sum_{i=1}^n P(\text{flower } i \text{ is selected and it is the biggest}) \quad (1)$$

$$= \sum_{i=1}^n P(\text{flower } i \text{ is selected} | i \text{ is the best}) P(\text{flower } i \text{ is the best}) \quad (2)$$

$$= \left(\sum_{i=1}^{r-1} P(\text{flower } i \text{ is selected} | i \text{ is the best}) + \sum_{i=r}^n P(\text{flower } i \text{ is selected} | i \text{ is the best}) \right) P(\text{flower } i \text{ is the best}) \quad (3)$$

Following this, can you derive the probability formulation for the probability that the biggest flower is actually selected.

Hint: Ethan will stop at any flower F larger than the calibration one. What's the probability of this one (F) to be the actual biggest one? Maybe all the flowers before F is smaller than M .

- (ii) [4 pts] Now compute the probability when n goes to ∞ .

Hint: you might need to do a bit of calculus here. For example, when $n \rightarrow \infty$, $\frac{1}{n}$ can be dt in an integral. Here t is a new variable we introduce.

- (iii) [3 pts] Take the derivative and tell Ethan when should he stop calibration. (what is r ?) Please show all steps of your work

- (b) (i) [1 pt] If it is a sunny day (+s), then Ethan will be happy (+h) with a high probability. If Ethan is happy (+h), he will execute his plan to pick a flower with a high probability (+p). Otherwise, he will execute his plan with a low probability (+p'). Sunny day (S), Ethan's happiness (H), and pick a flower (P) are three variables. Draw out the bayes net with those three variables.
- (ii) [2 pts] If condition on Ethan's happiness, is the weather and flower picking independent? Explain your answer intuitively **without referring to the active/inactive triples**.
- (iii) [7 pts] Now we find more concrete information about this activity. For weather, it can be sunny and rainy. For happiness, it can be happy and unhappy. For flower pick, it can be pick or not pick. $Pr(sunny) = 0.4$, $Pr(happy|sunny) = 0.8$, $Pr(happy|rainy) = 0.5$, $Pr(pick|happy) = 0.8$, $Pr(pick|unhappy) = 0.2$. Can you compute $Pr(sunny|pick)$? Please show all steps of your work.

- (c) Consider a modified setting where Ethan and his friends Fiona are **competing to pick a bigger flower**. Ethan and Fiona are walking in two separate fields with flowers of the same set of sizes $[s_1, s_2, s_3, s_4, \dots, s_n]$, which in turn have ranks $[1, 2, 3, 4, \dots, n]$ (where higher ranks means bigger sizes), but arranged in (possibly) different order. At each discrete time-step t , both Ethan and Fiona will see one more flower in their respective field.

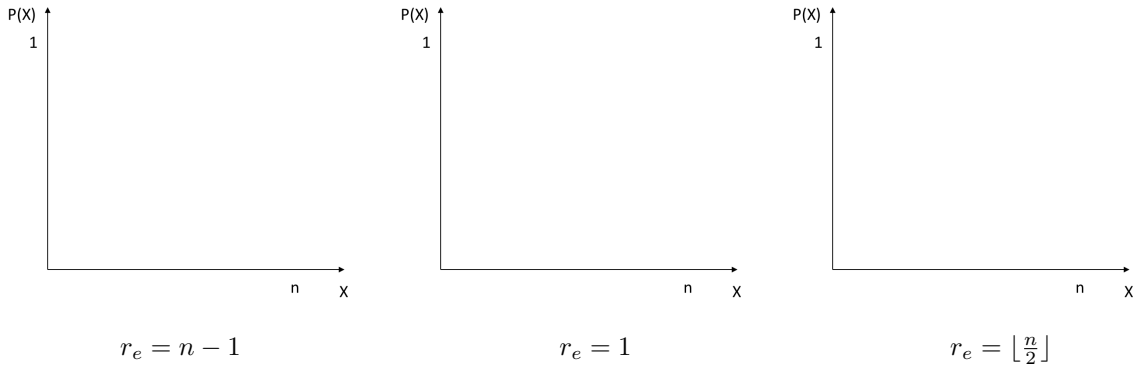
Ethan and Fiona cannot observe the other person's field directly (and thus cannot see what the other person has seen so far).

If Fiona knows that Ethan is committed to the strategy described in part (a.i), and the number of flowers Ethan is checking out is exactly r_e . What is Fiona's optimal r_f to get a bigger flower than Ethan?

- (i) [1 pt] Consider the rank of the flower Ethan ended up picking as random variable X , and the rank of the flower Fiona ended up picking as random variable Y . What is the condition for Fiona strictly beating Ethan in the competition, in term of X and Y

- (ii) [3 pts] How does the probability distribution for $P(X = x)$ for $x = 1, 2, 3, \dots, n$ depend on r_e ? Please provide a sketch for $P(X)$ parametrized by $r_e = n - 1$, $r_e = 1$, $r_e = \lfloor \frac{n}{2} \rfloor$. Note that the function will be discrete.

You do not have to derive the expression for $P(X = x)$ in terms of r_e . Reasoning the shape should be adequate



Justification for the shape of $r_e = n - 1$:

Justification for the shape of $r_e = 1$:

Justification for the shape of $r_e = \lfloor \frac{n}{2} \rfloor$:

- (iii) [2 pts] Does r_f parametrize $P(Y)$ the same way r_e parametrizes $P(X)$?

☐ Yes ☐ No

Explanation:

- (iv) [2 pts] What is the probability of $P(Y > X)$? Your answer should include two summations, and you are not required to simplify $P(X)$, $P(Y)$.

$$P(Y > X) =$$

- (v) [2 pts] Suppose you have the full expression for $P(X)$ and $P(Y)$ in terms of r_e , r_f respectively. How would you find the optimal r_f that would maximize the chance for Fiona to beat Ethan in this competition?

Explanation: