

Name: \_\_\_\_\_

Student ID: \_\_\_\_\_

## (0pts) 1. (Numerical Answer Formatting)

Many of the questions in this homework have answers that are decimal numbers. Due to current limitations of Gradescope, your answers must be an exact string match to ours. In order to ensure an exact match, please carefully follow the following formatting for your numerical answers.

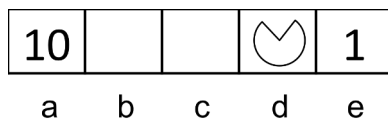
- Do not round decimals. None of the answers are infinite decimals, so include full precision (all answers should be less than 5 places after the decimal).
- Do not include any leading or trailing 0s unless they are necessary to show the location of the decimal.
- If the number is an integer, do not include a decimal

**Examples.** .1234, -.001, 10.4, -10, 0

**Note.** If you use the Python interpreter to do your math, floating point error may lead to inexact decimal numbers. It is probably best to use another calculator, but if you do use Python you may need to adjust its output to get the actual exact answer.

## (9pts) 2. (Solving MDPs)

Consider the gridworld MDP for which the actions **left** and **right** are always successful. Specifically, the available actions in each state are to move to the neighboring grid squares. From state  $a$ , there is also an **exit** action available, which results in going to the terminal state and collecting a reward of 10. Similarly, in state  $e$ , the reward for the **exit** action is 1. (Exit actions are always successful.)



Let the discount factor be  $\gamma = 1$ . Fill in the following quantities.

$V_0(d) =$  \_\_\_\_\_

$V_1(d) =$  \_\_\_\_\_

$V_2(d) =$  \_\_\_\_\_

$V_3(d) =$  \_\_\_\_\_

$V_4(d) =$  \_\_\_\_\_

$V_5(d) =$  \_\_\_\_\_

(9pts) 3. (Value Iteration Convergence Values)

Consider the gridworld where **left** and **right** actions are always successful. Specifically, the available actions in each state are to move to the neighboring grid squares. From state  $a$ , there is also an **exit** action available, which results in going to the terminal state and collecting a reward of 10. Similarly, in state  $e$ , the reward for the **exit** action is 1. (Exit actions are always successful.)

10				1
a	b	c	d	e

Let the discount factor  $\gamma = 0.2$ . Fill in the following quantities.

$$V^*(a) = V_\infty(a) = \underline{\hspace{2cm}}$$

$$V^*(b) = V_\infty(b) = \underline{\hspace{2cm}}$$

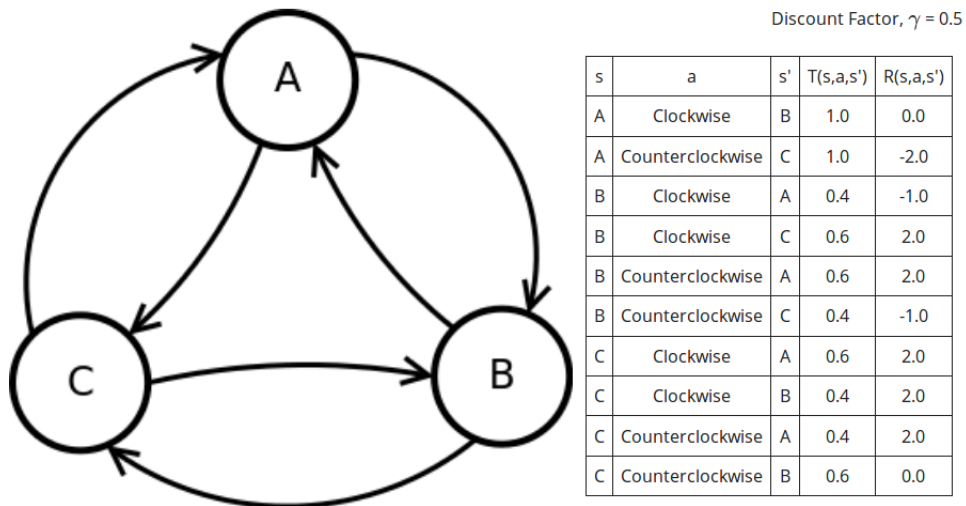
$$V^*(c) = V_\infty(c) = \underline{\hspace{2cm}}$$

$$V^*(d) = V_\infty(d) = \underline{\hspace{2cm}}$$

$$V^*(e) = V_\infty(e) = \underline{\hspace{2cm}}$$

(12pts) 4. (Value Iteration: Cycle)

Consider the following transition diagram, transition function and reward function for an MDP.



Suppose after iteration  $k$  of value iteration we end up with the following values for  $V_k$ :

$V_k(A)$	$V_k(B)$	$V_k(C)$
0.400	1.400	2.160

(a) What is  $V_{k+1}(A)$ ? *Answer.* \_\_\_\_\_

(b) Now, suppose that we ran value iteration to completion and found the following value function,  $V^*$ .

$V^*(A)$	$V^*(B)$	$V^*(C)$
0.881	1.761	2.616

What is  $Q^*(A, \text{clockwise})$ ?

*Answer.* \_\_\_\_\_

(c) What is  $Q^*(A, \text{counterclockwise})$ ?

*Answer.* \_\_\_\_\_

(d) What is the optimal action from state A?

☐ Clockwise    ☐ Counterclockwise

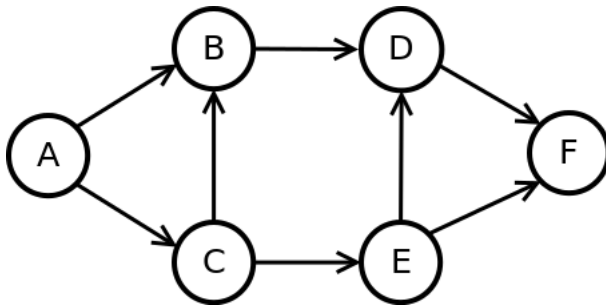
(7pts) 5. (Value Iterations Properties)

Which of the following are true about value iteration? We assume the MDP has a finite number of actions and states, and that the discount factor satisfies  $0 < \gamma < 1$ .

- ☐ Value iteration is guaranteed to converge.
- ☐ Value iteration will converge to the same vector of values ( $V^*$ ) no matter what values we use to initialize  $V$ .
- ☐ None of the above

(6pts) 6. (Value Iteration Convergence)

We will consider a simple MDP that has six states,  $A$ ,  $B$ ,  $C$ ,  $D$ ,  $E$ , and  $F$ . Each state has a single action, **go**. An arrow from a state  $x$  to a state  $y$  indicates that it is possible to transition from state  $x$  to next state  $y$  when **go** is taken. If there are multiple arrows leaving a state  $x$ , transitioning to each of the next states is equally likely. The state  $F$  has no outgoing arrows: once you arrive in  $F$ , you stay in  $F$  for all future times. The reward is one for all transitions, with one exception: staying in  $F$  gets a reward of zero. Assume a discount factor  $= 0.5$ . We assume that we initialize the value of each state to 0. (Note: you should not need to explicitly run value iteration to solve this problem.)



- (a) After how many iterations of value iteration will the value for state  $E$  have become exactly equal to the true optimum? (Enter inf if the values will never become equal to the true optimal but only converge to the true optimal.)

Answer. \_\_\_\_\_

- (b) How many iterations of value iteration will it take for the values of all states to converge to the true optimal values? (Enter inf if the values will never become equal to the true optimal but only converge to the true optimal.)

Answer. \_\_\_\_\_

(10pts) 7. (Policy Evaluation)

Consider the gridworld where Left and Right actions are always successful.

Specifically, the available actions in each state are to move to the neighboring grid squares. From state  $a$ , there is also an exit action available, which results in going to the terminal state and collecting a reward of 10. Similarly, in state  $e$ , the reward for the exit action is 1. (Exit actions are always successful.)

The discount factor is  $\gamma = 1$ .

10				1
a	b	c	d	e

- (a) Consider the policy  $\pi_1$  shown below, and evaluate the following quantities for this policy.

→	→	→	→	exit
a	b	c	d	e

$$V^{\pi_1}(a) = \underline{\hspace{2cm}}$$

$$V^{\pi_1}(b) = \underline{\hspace{2cm}}$$

$$V^{\pi_1}(c) = \underline{\hspace{2cm}}$$

$$V^{\pi_1}(d) = \underline{\hspace{2cm}}$$

$$V^{\pi_1}(e) = \underline{\hspace{2cm}}$$

- (b) Consider the policy  $\pi_2$  shown below, and evaluate the following quantities for this policy.

exit	←	←	→	exit
a	b	c	d	e

$$V^{\pi_2}(a) = \underline{\hspace{2cm}}$$

$$V^{\pi_2}(b) = \underline{\hspace{2cm}}$$

$$V^{\pi_2}(c) = \underline{\hspace{2cm}}$$

$$V^{\pi_2}(d) = \underline{\hspace{2cm}}$$

$$V^{\pi_2}(e) = \underline{\hspace{2cm}}$$

(9pts) 8. (Policy Iteration)

Consider the gridworld where **left** and **right** actions are always successful.

Specifically, the available actions in each state are to move to the neighboring grid squares. From state  $a$ , there is also an **exit** action available, which results in going to the terminal state and collecting a reward of 10. Similarly, in state  $e$ , the reward for the **exit** action is 1. (Exit actions are always successful.)

The discount factor is  $\gamma = 0.9$ .

10				1
a	b	c	d	e

We will execute one round of policy iteration.

Consider the policy  $\pi_i$  shown below, and evaluate the following quantities for this policy.

exit	←	→	←	exit
a	b	c	d	e

$$V^{\pi_i}(a) = \underline{\hspace{2cm}}$$

$$V^{\pi_i}(b) = \underline{\hspace{2cm}}$$

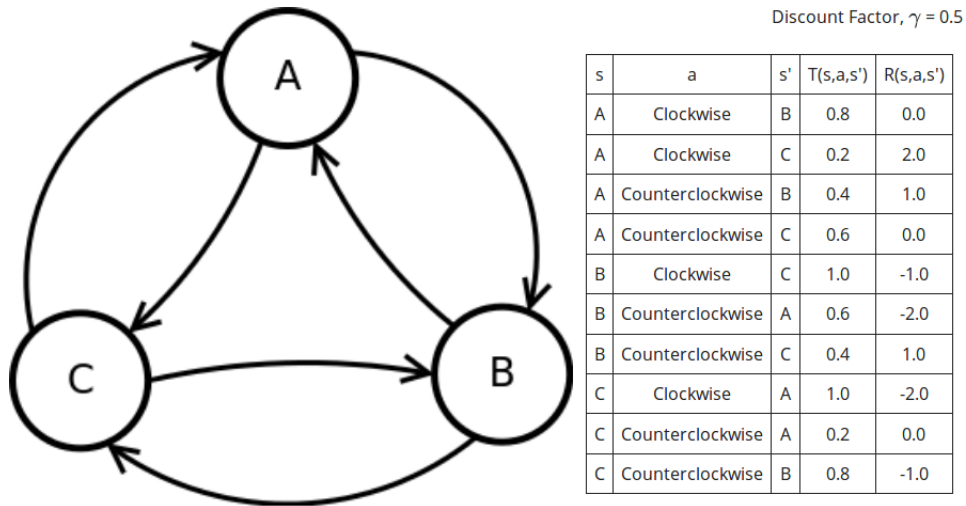
$$V^{\pi_i}(c) = \underline{\hspace{2cm}}$$

$$V^{\pi_i}(d) = \underline{\hspace{2cm}}$$

$$V^{\pi_i}(e) = \underline{\hspace{2cm}}$$

(14pts) 9. (Policy Iteration: Cycle)

Consider the following transition diagram, transition function and reward function for an MDP.



Suppose we are doing policy evaluation, by following the policy given by the left-hand side table below. Our current estimates (at the end of some iteration of policy evaluation) of the value of states when following the current policy is given in the right-hand side table.

A	B	C
Counterclockwise	Counterclockwise	Counterclockwise

$V_k^\pi(A)$	$V_k^\pi(B)$	$V_k^\pi(C)$
0.000	-0.840	-1.080

(a) What is  $V_{k+1}^\pi(A)$ ? *Answer.* \_\_\_\_\_

(b) Suppose that policy evaluation converges to the following value function,  $V_\infty^\pi$ .

$V_\infty^\pi(A)$	$V_\infty^\pi(B)$	$V_\infty^\pi(C)$
-0.203	-1.114	-1.266

Now let's execute policy improvement.

What is  $Q_\infty^\pi(A, \text{clockwise})$ ? *Answer.* \_\_\_\_\_

(c) What is  $Q_\infty^\pi(A, \text{counterclockwise})$ ? *Answer.* \_\_\_\_\_

(d) What is the updated action for state A? ☐ Clockwise ☐ Counterclockwise

(7pts) 10. (Wrong Discount Factor)

Bob notices value iteration converges more quickly with smaller  $\gamma$  and rather than using the true discount factor  $\gamma$ , he decides to use a discount factor of  $\alpha \gamma$  with  $0 < \alpha < 1$  when running value iteration. Mark each of the following that are guaranteed to be true:

- ☐ While Bob will not find the optimal value function, he could simply rescale the values he finds by  $\frac{1-\gamma}{1-\alpha}$  to find the optimal value function.
- ☐ If the MDP's transition model is deterministic and the MDP has zero rewards everywhere, except for a single transition at the goal with a positive reward, then Bob will still find the optimal policy.
- ☐ If the MDP's transition model is deterministic, then Bob will still find the optimal policy.
- ☐ Bob's policy will tend to more heavily favor short-term rewards over long-term rewards compared to the optimal policy.
- ☐ None of the above.

(10pts) 11. (MDP Properties)

(a) Which of the following statements are true for an MDP?

- ☐ If the only difference between two MDPs is the value of the discount factor then they must have the same optimal policy.
- ☐ For an infinite horizon MDP with a finite number of states and actions and with a discount factor  $\gamma$  that satisfies  $0 < \gamma < 1$ , value iteration is guaranteed to converge.
- ☐ When running value iteration, if the policy (the greedy policy with respect to the values) has converged, the values must have converged as well.
- ☐ None of the above

(b) Which of the following statements are true for an MDP?

- ☐ If one is using value iteration and the values have converged, the policy must have converged as well.
- ☐ Expectimax will generally run in the same amount of time as value iteration on a given MDP.
- ☐ For an infinite horizon MDP with a finite number of states and actions and with a discount factor  $\gamma$  that satisfies  $0 < \gamma < 1$ , policy iteration is guaranteed to converge.
- ☐ None of the above



(7pts) 12. John, James, Alvin and Michael all get to act in an MDP  $(S, A, T, \gamma, R, s_0)$ .

**John** runs value iteration until he finds  $V^*$  which satisfies

$$\forall s (s \in S \rightarrow V^*(s) = \max_{a \in A} \sum_{s'} T(s, a, s') (R(s, a, s') + \gamma V^*(s')))$$

and acts according to

$$\pi_{\text{John}} = \arg \max_{a \in A} \sum_{s'} T(s, a, s') (R(s, a, s') + \gamma V^*(s')).$$

**James** acts according to an arbitrary policy  $\pi_{\text{James}}$ .

**Alvin** takes James's policy  $\pi_{\text{James}}$  and runs one round of policy iteration to find his policy  $\pi_{\text{Alvin}}$ .

**Michael** takes John's policy and runs one round of policy iteration to find his policy  $\pi_{\text{Michael}}$ .

Note: One round of policy iteration = performing policy evaluation followed by performing policy improvement.

Mark all of the following that are guaranteed to be true:

- ☐ It is guaranteed that  $\forall s \in S : V^{\pi_{\text{James}}}(s) \geq V^{\pi_{\text{Alvin}}}(s)$
- ☐ It is guaranteed that  $\forall s \in S : V^{\pi_{\text{Michael}}}(s) \geq V^{\pi_{\text{Alvin}}}(s)$
- ☐ It is guaranteed that  $\forall s \in S : V^{\pi_{\text{Michael}}}(s) > V^{\pi_{\text{John}}}(s)$
- ☐ It is guaranteed that  $\forall s \in S : V^{\pi_{\text{James}}}(s) > V^{\pi_{\text{John}}}(s)$
- ☐ None of the above.