

## Q1. MDP

Pacman is using MDPs to maximize his expected utility. In each environment:

- Pacman has the standard actions {North, East, South, West} unless blocked by an outer wall
  - There is a reward of 1 point when eating the dot (for example, in the grid below,  $R(C, South, F) = 1$ )
  - The game ends when the dot is eaten
- (a) Consider the following grid where there is a single food pellet in the bottom right corner ( $F$ ). The **discount** factor is 0.5. There is no living reward. The states are simply the grid locations.

A	B	C
D	E	F ○

(i) What is the optimal policy for each state?

State	$\pi(state)$
A	
B	
C	
D	
E	

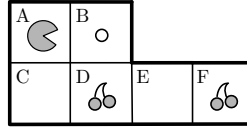
(ii) What is the optimal value for the state of being in the upper left corner ( $A$ )? Reminder: the discount factor is 0.5.

$$V^*(A) =$$

(iii) Using value iteration with the value of all states equal to zero at  $k=0$ , for which iteration  $k$  will  $V_k(A) = V^*(A)$ ?

$$k =$$

- (b) Consider a new Pacman level that begins with cherries in locations  $D$  and  $F$ . Landing on a grid position with cherries is worth 5 points and then the cherries at that position disappear. There is still one dot, worth 1 point. The game still only ends when the dot is eaten.



- (i) With no discount ( $\gamma = 1$ ) and a living reward of -1, what is the optimal policy for the states in this level's state space?

- (ii) With no discount ( $\gamma = 1$ ), what is the range of living reward values such that Pacman eats exactly one cherry when starting at position  $A$ ?

- (c) Quick reinforcement learning questions [PLEASE WRITE CLEARLY]:

- (i) What is the difference between value-iteration and TD-learning?

- (ii) What is the difference between TD-learning and Q-learning?

- (iii) What is the purpose of using a learning rate ( $\alpha$ ) during Q-learning?

- (iv) In value iteration, we store the value of each state. What do we store during *approximate* Q-learning?

- (v) Give one advantage and one disadvantage of using approximate Q-learning rather than standard Q-learning.

## Q2. Markov Decision Processes

Consider a simple MDP with two states,  $S_1$  and  $S_2$ , two actions,  $A$  and  $B$ , a discount factor  $\gamma$  of  $1/2$ , reward function  $R$  given by

$$R(s, a, s') = \begin{cases} 1 & \text{if } s' = S_1; \\ -1 & \text{if } s' = S_2; \end{cases}$$

and a transition function specified by the following table.

$s$	$a$	$s'$	$T(s, a, s')$
$S_1$	$A$	$S_1$	$1/2$
$S_1$	$A$	$S_2$	$1/2$
$S_1$	$B$	$S_1$	$2/3$
$S_1$	$B$	$S_2$	$1/3$
$S_2$	$A$	$S_1$	$1/2$
$S_2$	$A$	$S_2$	$1/2$
$S_2$	$B$	$S_1$	$1/3$
$S_2$	$B$	$S_2$	$2/3$

- (a) Perform a single iteration of value iteration, filling in the resultant Q-values and state values in the following tables. Use the specified initial value function  $V_0$ , rather than starting from all zero state values. Only compute the entries not labeled “skip”.

$s$	$a$	$Q_1(s, a)$
$S_1$	$A$	
$S_1$	$B$	
$S_2$	$A$	skip
$S_2$	$B$	skip

$s$	$V_0(s)$	$V_1(s)$
$S_1$	2	
$S_2$	3	skip

- (b) Suppose that Q-learning with a learning rate  $\alpha$  of  $1/2$  is being run, and the following episode is observed.

$s_1$	$a_1$	$r_1$	$s_2$	$a_2$	$r_2$	$s_3$
$S_1$	$A$	1	$S_1$	$A$	-1	$S_2$

Using the initial Q-values  $Q_0$ , fill in the following table to indicate the resultant progression of Q-values.

$s$	$a$	$Q_0(s, a)$	$Q_1(s, a)$	$Q_2(s, a)$
$S_1$	$A$	$-1/2$		
$S_1$	$B$	0		
$S_2$	$A$	-1		
$S_2$	$B$	1		

- (c) Given an arbitrary MDP with state set  $S$ , transition function  $T(s, a, s')$ , discount factor  $\gamma$ , and reward function  $R(s, a, s')$ , and given a constant  $\beta > 0$ , consider a modified MDP  $(S, T, \gamma, R')$  with reward function  $R'(s, a, s') = \beta \cdot R(s, a, s')$ . Prove that the modified MDP  $(S, T, \gamma, R')$  has the same set of optimal policies as the original MDP  $(S, T, \gamma, R)$ .

- (d) Although in this class we have defined MDPs as having a reward function  $R(s, a, s')$  that can depend on the initial state  $s$  and the action  $a$  in addition to the destination state  $s'$ , MDPs are sometimes defined as having a reward function  $R(s')$  that depends only on the destination state  $s'$ . Given an arbitrary MDP with state set  $S$ , transition function  $T(s, a, s')$ , discount factor  $\gamma$ , and reward function  $R(s, a, s')$  that *does depend* on the initial state  $s$  and the action  $a$ , define an *equivalent* MDP with state set  $S'$ , transition function  $T'(s, a, s')$ , discount factor  $\gamma'$ , and reward function  $R'(s')$  that depends only on the destination state  $s'$ .

By *equivalent*, it is meant that there should be a one-to-one mapping between state-action sequences in the original MDP and state-action sequences in the modified MDP (with the same value). **You do not need to give a proof of the equivalence.**

**States:**  $S' =$

**Transition function:**  $T'(s, a, s') =$

**Discount factor:**  $\gamma' =$

**Reward function:**  $R'(s') =$