## Q1. Q-uagmire

Consider an unknown MDP with three states (A, B and C) and two actions  $(\leftarrow \text{ and } \rightarrow)$ . Suppose the agent chooses actions according to some policy  $\pi$  in the unknown MDP, collecting a dataset consisting of samples (s, a, s', r) representing taking action a in state s resulting in a transition to state s' and a reward of r.

S	a	s'	r
$\overline{A}$	$\rightarrow$	В	2
$\boldsymbol{C}$	$\leftarrow$	$\boldsymbol{B}$	2
$\boldsymbol{B}$	$\rightarrow$	$\boldsymbol{C}$	-2
$\boldsymbol{A}$	$\rightarrow$	$\boldsymbol{B}$	4

You may assume a discount factor of  $\gamma = 1$ .

(a) Recall the update function of Q-learning is:

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left(r_t + \gamma \max_{a'} Q(s_{t+1}, a')\right)$$

Assume that all Q-values are initialized to 0, and use a learning rate of  $\alpha = \frac{1}{2}$ .

(i) Run Q-learning on the above experience table and fill in the following Q-values:

$$Q(A, \to) = \underbrace{\qquad \qquad 5/2 \qquad \qquad Q(B, \to) = \underbrace{\qquad \qquad -1/2}}_{Q_1(A, \to) = \frac{1}{2} \cdot Q_0(A, \to) + \frac{1}{2} \left( 2 + \gamma \max_{a'} Q(B, a') \right) = 1}_{Q_1(C, \leftarrow) = 1}$$

$$Q_1(B, \to) = \frac{1}{2} (-2 + 1) = -\frac{1}{2}$$

$$Q_2(A, \to) = \frac{1}{2} \cdot 1 + \frac{1}{2} \left( 4 + \max_{a'} Q_1(B, a') \right)$$

$$= \frac{1}{2} + \frac{1}{2} (4 + 0) = \frac{5}{2}.$$

(ii) After running Q-learning and producing the above Q-values, you construct a policy  $\pi_Q$  that maximizes the Q-value in a given state:

$$\pi_Q(s) = \arg\max_a Q(s, a).$$

What are the actions chosen by the policy in states A and B?

 $\pi_Q(A)$  is equal to:

 $\bigcirc \ \pi_O(A) = \leftarrow.$ 

 $\bigcirc \quad \pi_O(A) = \text{Undefined.}$ 

- **(b)** This question considers properties of reinforcement learning algorithms for *arbitrary* discrete MDPs; you do not need to refer to the MDP considered in the previous parts.
  - (i) Which of the following methods, at convergence, provide enough information to obtain an optimal policy? (Assume adequate exploration.)

Model-based learning of T(s, a, s') and R(s, a, s').

 $\square$  Direct Evaluation to estimate V(s).

 $\square$  Temporal Difference learning to estimate V(s).

- Q-Learning to estimate Q(s, a). Given enough data, model-based learning will get arbitrarily close to the true model of the environment, at which point planning (e.g. value iteration) can be used to find an optimal policy. Q-learning is similarly guaranteed to converge to the optimal Q-values of the optimal policy, at which point the optimal policy can be recovered by  $\pi^*(s) = \arg\max_a Q(s, a)$ . Direct evaluation and temporal difference learning both only recover a value function V(s), which is insufficient to choose between actions without knowledge of the transition probabilities.
- (ii) In the limit of infinite timesteps, under which of the following exploration policies is Q-learning guaranteed to converge to the optimal Q-values for all state? (You may assume the learning rate  $\alpha$  is chosen appropriately, and that the MDP is ergodic: i.e., every state is reachable from every other state with non-zero probability.)

A fixed policy taking actions uniformly at random.

A greedy policy.

An  $\epsilon$ -greedy policy

A fixed optimal policy. For Q-learning to converge, every state-action pair (s, a) must occur infinitely often. A uniform random policy will achieve this in an ergodic MDP. A fixed optimal policy will not take any suboptimal actions and so will not explore enough. Similarly a greedy policy will stop taking actions the current Q-values suggest are suboptimal, and so will never update the Q-values for supposedly suboptimal actions. (This is problematic if, for example, an action most of the time yields no reward but occasionally yields very high reward. After observing no reward a few times, Q-learning with a greedy policy would stop taking that action, never obtaining the high reward needed to update it to its true value.)

## Q2. RL: Amusement Park

After the disastrous waterslide experience you decide to go to an amusement park instead. In the previous questions the MDP was based on a single ride (a water slide). Here our MDP is about choosing a ride from a set of many rides.

You start off feeling well, getting positive rewards from rides, some larger than others. However, there is some chance of each ride making you sick. If you continue going on rides while sick there is some chance of becoming well again, but you don't enjoy the rides as much, receiving lower rewards (possibly negative).

You have never been to an amusement park before, so you don't know how much reward you will get from each ride (while well or sick). You also don't know how likely you are to get sick on each ride, or how likely you are to become well again. What you do know about the rides is:

Actions / Rides	Type	Wait	Speed
Big Dipper	Rollercoaster	Long	Fast
Wild Mouse	Rollercoaster	Short	Slow
Hair Raiser	Drop tower	Short	Fast
Moon Ranger	Pendulum	Short	Slow
Leave the Park	Leave	Short	Slow

We will formulate this as an MDP with two states, well and sick. Each ride corresponds to an action. The 'Leave the Park' action ends the current run through the MDP. Taking a ride will lead back to the same state with some probability or take you to the other state. We will use a feature based approximation to the Q-values, defined by the following four features and associated weights:

Features	Initial Weights
$f_0(state, action) = 1$ (this is a bias feature that is always 1)	$w_0 = 1$
$f_1(state, action) = \begin{cases} 1 & \text{if } action \text{ type is Rollercoaster} \\ 0 & \text{otherwise} \end{cases}$	$w_1 = 2$
$f_2(state, action) = \begin{cases} 1 & \text{if } action \text{ wait is Short} \\ 0 & \text{otherwise} \end{cases}$	$w_2 = 1$
$f_3(state, action) = \begin{cases} 1 & \text{if } action \text{ speed is Fast} \\ 0 & \text{otherwise} \end{cases}$	$w_3 = 0.5$

(a) Calculate Q('Well', 'Big Dipper'):

$$1 + 2 + 0 + 0.5 = 3.5$$

(b) Apply a Q-learning update based on the sample ('Well', 'Big Dipper', 'Sick', -10.5), using a learning rate of  $\alpha = 0.5$  and discount of  $\gamma = 0.5$ . What are the new weights?

```
Difference = -10.5 + 0.5 * max(4, 3.5, 2.5, 2.0, 2.0) - 3.5 = -12

w_0 = 1 - 6 * 1 = -5

w_1 = 2 - 6 * 1 = -4

w_2 = 1 - 6 * 0 = 1

w_3 = 0.5 - 6 * 1 = -5.5
```

