

1. (20%) The following algorithm structure presents one-step temporal-difference control methods. Name all possible algorithms (Sarsa, Q-learning, or on-policy Expected Sarsa) that can fit this structure and describe the corresponding update rule(s).

```

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ 
Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$ 
Loop for each episode:
  Initialize  $S$ 
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
  Loop for each step of episode:
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    
 $Q(S, A) \leftarrow Q(S, A) + \alpha (R + V(S') - Q(S, A))$ 

     $S \leftarrow S'; A \leftarrow A';$ 
  until  $S$  is terminal

```

2. (20%) Name the following algorithm:

```

Input: a policy  $\pi$  to be evaluated
Initialize:
   $V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$ 
   $Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$ 
Loop forever (for each episode):
  Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
   $G \leftarrow 0$ 
  Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
     $G \leftarrow G + R_{t+1}$ 
    Unless  $S_t$  appears in  $S_0, S_1, \dots, S_{t-1}$ :
      Append  $G$  to  $Returns(S_t)$ 
       $V(S_t) \leftarrow \text{average}(Returns(S_t))$ 

```

3.
 - (a) (10%) Let p denote the state transition probability. Given a starting state s_t and a target policy π , derive the probability of having the state-action trajectory $a_t, s_{t+1}, a_{t+1}, \dots, s_T$.
 - (b) (10%) Use (a) to show that the importance-sampling ratio (i.e., the relative probability of the trajectory under the target policy π and behavior policy b) is

$$\rho_{t:T-1} = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{b(A_k|S_k)}.$$

4. (20%) The following algorithm is called double Q-learning that has been widely used to address the problem of maximization bias. Fill out the blanks to complete the algorithm

```

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ 
Initialize  $Q_1(s, a)$  and  $Q_2(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q_1(\text{terminal}, \cdot) = 0$ 
Loop for each episode:
    Initialize  $S$ 
    Loop for each step of episode:
        Choose  $A$  from  $S$  using the policy  $\varepsilon$ -greedy in  $Q_1 + Q_2$ 
        Take action  $A$ , observe  $R, S'$ 
        With 0.5 probability:
             $Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left( R + \gamma \max_{a \in \mathcal{A}(S')} Q_1(S', a) - Q_1(S, A) \right)$ 
        else:
             $Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left( R + \gamma \max_{a \in \mathcal{A}(S')} Q_2(S', a) - Q_2(S, A) \right)$ 
         $S \leftarrow S'$ 
    until  $S$  is terminal

```

5. (20%) The following algorithm structure presents an n-step bootstrapping control methods. Name all possible algorithms (on-policy Sarsa, off-policy Sarsa, or on-policy Expected Sarsa) that can fit this structure and describe the corresponding update rules (a pair of blanks for one particular algorithm).

```

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}$ 
Initialize  $\pi$  to be  $\varepsilon$ -greedy with respect to  $Q$ , or to a fixed given policy
Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ , a positive integer  $n$ 
All store and access operations (for  $S_t$ ,  $A_t$ , and  $R_t$ ) can take their index mod  $n$ 

Loop for each episode:
    Initialize and store  $S_0 \neq \text{terminal}$ 
    Select and store an action  $A_0 \sim \pi(\cdot | S_0)$ 
     $T \leftarrow \infty$ 
    Loop for each step of episode,  $t = 0, 1, 2, \dots$ :
        | If  $t < T$ , then:
        |     Take action  $A_t$ 
        |     Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$ 
        |     If  $S_{t+1}$  is terminal, then:
        |          $T \leftarrow t + 1$ 
        |     else:
        |         Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ 
         $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
        | If  $\tau \geq 0$ :
        |      $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
        |     If  $\tau + n < T$ , then  $G \leftarrow G + \text{[ ]}$  ( $G_{\tau:\tau+n}$ )
        |      $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha \text{[ ]}$ 
        |     If  $\pi$  is being learned, then ensure that  $\pi(\cdot | S_\tau)$  is  $\varepsilon$ -greedy wrt  $Q$ 
    Until  $\tau = T - 1$ 

```