#### 535514 Spring 2024: Reinforcement Learning

(Due: 2024/03/20, Wednesday, 21:00)

Homework 1: Warm-Up – Fundamentals of MDPs and RL

Submission Guidelines: Your deliverables shall consist of 2 separate files – (i) A PDF file: Please compile all your write-ups into one .pdf file (photos/scanned copies are acceptable; please make sure that the electronic files are of good quality and reader-friendly); (ii) A zip file: Please compress all your source code into one .zip file. Please submit your deliverables via E3.

### Problem 1 (Q-Value Iteration)

(15+15=30 points)

(a) Recall that in Lecture 3, we define  $V^*(s) := \max_{\pi} V^{\pi}(s)$  and  $Q^*(s, a) := \max_{\pi} Q^{\pi}(s, a)$ . Suppose  $\gamma \in (0, 1)$ . Prove the following Bellman optimality equations:

$$V^*(s) = \max_{a} Q^*(s, a) \tag{1}$$

$$Q^*(s,a) = R_{s,a} + \gamma \sum_{s'} P^a_{ss'} V^*(s').$$
(2)

Please carefully justify every step of your proof. (Hint: For (1), you may first prove that  $V^*(s) \leq \max_a Q^*(s, a)$  and then show  $V^*(s) < \max_a Q^*(s, a)$  cannot happen by contradiction. On the other hand, (2) can be shown by using the similar argument or by leveraging the fact that  $Q^{\pi}(s, a) = R_{s,a} + \gamma \sum_{s'} P_{ss'}^a V^{\pi}(s')$ 

Based on (a), we thereby have the recursive Bellman optimality equation for the optimal action-value function  $Q_*$  as:

$$Q^*(s,a) = R_{s,a} + \gamma \sum_{s'} P^a_{ss'} \Big( \max_{a'} Q^*(s',a') \Big)$$
(3)

Similar to the standard Value Iteration, we can also study the Q-Value Iteration by defining the Bellman optimality operator  $T^*: \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} \to \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  for the action-value function: for every state-action pair (s, a)

$$[T^*(Q)](s,a) := R_{s,a} + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q(s',a')$$
(4)

Show that the operator  $T^*$  is a  $\gamma$ -contraction operator in terms of  $\infty$ -norm. Please carefully justify every step of your proof. (Hint: For any two action-value functions Q, Q', we have  $\|T^*(Q) - T^*(Q')\|_{\infty} = \max_{(s,a)} \left| [T^*(Q)](s,a) - [T^*(Q')](s,a) \right|$ )

## Problem 2 (Regularized MDPs)

(10+10=20 points)

In Lecture 4, we formally describe the regularized MDP, which is a direct extension of the classic MDP with a regularizer  $\Omega$ . In this problem, for simplicity, suppose we use the Shannon entropy as our regularizer, i.e.,  $\Omega(\pi(\cdot|s)) \equiv H(\pi(\cdot|s)) := -\sum_{a \in \mathcal{A}} \pi(a|s) \ln \pi(a|s)$ . Let us verify a few important properties mentioned in Lecture 4 as follows.

(2) Recall that we introduce the "regularized Bellman expectation operator"  $T_{\Omega}^{\pi}$  as

$$[T_{\Omega}^{\pi}V](s) := R_s^{\pi} + \Omega(\pi(\cdot|s)) + \gamma P_{ss'}^{\pi}V. \tag{5}$$

Please verify that  $T_{\Omega}^{\pi}$  is a contraction operator in  $L_{\infty}$  norm. (Hint: Try to extend the proof procedure of the contraction property of  $T^{\pi}$  in Lecture 3)

Moreover, under regularized MDPs, we study the optimal value functions  $V_{\Omega}^*$  and optimal Q functions  $Q_{\Omega}^*$ 

and learn the Bellman optimality equations as

$$V_{\Omega}^*(s) = \max_{\pi \in \Pi} R_s^{\pi} + \gamma P_s^{\pi} V_{\Omega}^* \tag{6}$$

$$Q_{\Omega}^{*}(s,a) = R_{s,a} + \gamma E_{s' \sim P(\cdot|s,a)}[V_{\Omega}^{*}(s')]. \tag{7}$$

Could you design an iterative algorithm that can obtain  $V_{\Omega}^*$  and  $Q_{\Omega}^*$ ? Please clearly write down the complete pseudo code of your algorithm and provide comments on each line of your pseudo code. (Hint: Try to extend the Value Iteration for standard MDPs to the regularized MDPs based on Equation 6)

#### Problem 3 (A Property Used in Policy Gradient)

(10 points)

Show the following useful property discussed in Lectures 5-6: for any function  $f: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ .

$$\mathbb{E}_{\tau \sim P_{\mu}^{\pi_{\theta}}} \left[ \sum_{t=0}^{\infty} \gamma^{t} f(s_{t}, a_{t}) \right] = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta}}} \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid s)} \left[ f(s, a) \right]$$
(8)

(Hint: It might be slightly easier to go from the RHS to LHS. Specifically, you may first expand the RHS of (8) into a sum of f(s, a) over s and a and then apply the definition of  $d^{\pi_{\theta}}_{\mu}$ , which involves a sum of probability over t. Next, try to reorganize the triple summation into the form of the LHS of (8))

#### Problem 4 (Implementing Policy Iteration and Value Iteration)

(35 points)

In this problem, we will implement Policy Iteration and Calue Iteration for a classic MDP environment called "Taxi" (Dietterich, 2000). This environment has been included in the OpenAI Gym: https://gym.openai.com/envs/Taxi-v3/. To accomplish this task, you may take the following steps:

Get familiar with the Taxi environment by reading the Gym documentation at https://www.gymlibrary.dev/environments/toy\_text/taxi/. The state space consists of 500 possible states as there are 25 taxi positions, 5 possible locations of the passenger (including the case when the passenger is in the taxi), and 4 destination locations. Moreover, the agent has 6 possible actions (namely, 0: move south; 1: move north; 2: move east; 3: move west; 4: pickup passenger; 5: drop off passenger). The rewards are: (i) -1 per step unless other reward is triggered; (ii) +20 for delivering passenger; (iii) -10 for executing "pickup" and "drop-off" actions illegally.

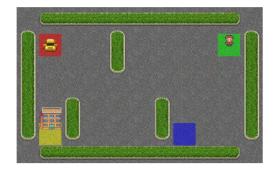


Figure 1: An illustration of the Taxi environment.

Read through **policy\_and\_value\_iteration**. **py** and then implement the two functions **policy\_iteration** and **value\_iteration** based on the pseudo code of PI and VI provided in the lecture slides.

Note: Please set  $\gamma = 0.9$  and the termination criterion  $\varepsilon = 10^{-3}$ . Moreover, you could use either Taxi-v2 or Taxi-v3 environment (Taxi-v3 is recommended). Note that discrepancy = 0 is a necessary condition (but not sufficient) of correct implementation, and with the default  $\varepsilon = 10^{-3}$ , you shall be able to observe zero discrepancy between the policies obtained by PI and VI.

#### Problem 5 (Installing and Getting Familiar With D4RL Dataset)

(10 points)

In this problem, you will be asked to install and investigate the D4RL Dataset, which is currently the standard dataset for Offline RL. To accomplish this task, you shall take the following steps:

- Please first take a careful look at the GitHub repo of D4RL https://github.com/Farama-Foundation/D4RL and the paper https://arxiv.org/abs/2004.07219. Then, install d4rl by following the installation guide on the GitHub repo (you may need to install several other packages, such as Gymnasium and MuJoCo).
  - Read and run the provided d4rl\_sanity\_check.py and finish the following tasks:
    - (1) Generate an offline dataset by directly running d4rl\_sanity\_check.py.
    - -1/(2) Describe the format of the dataset that you just obtained.
    - -1/3) Modify the **d4rl\_sanity\_check.py** and generate another offline dataset of one of the MuJoCo tasks (e.g., Hopper, Walker, or Halfcheetah). Similarly, describe the format of the MuJoCo dataset that you just obtained.
- One final remark: We will keep using the D4RL dataset for HW2, HW3, and the team final project. Therefore, it would be very helpful to now set up the environment that you could easily reuse subsequently.

```
(,(a) I Obv (, V"(s) < max, Q*(s,a)
   pf We have Vis = max (Vis) (by definition of Visi)
                    = max_(IT(a(s)Q(s,a))(by definition of V(s))
              max Q*(s, a)= max (max Q"(s, a) (by definition of Q"(s, a))
     Since T TS conditional distribution over all possible
     action given S, we must have
               0 = T(a(s) = 1 Va, Ys, ZT(a(s) = 1 VzeS
     => max = ( Tr(a(s)Q(s,a) < max = (max Q(s,a))
                                                      all Q(s,a)
            = max (max, Q'(s,a)/snce we lak

V*(s) (by 0) = max Q*(s,a) (by 0)
     =>V (S) & max Q (S, a) [
   ObZ. Visi maxa Qis,a) is impossible.
   pf By the previous discussion, that is,
      IT S.t. max (IT (als)Q(s,a)= max, (max, Q(s,a))
     Ne con let T((als): ) 1, if a: argmaxQ((s,u),
     the above equality 75 hold,
   By Obl and ObvZ, we know Vis=maxQis, a) -
```

(QI Dbl. Q'(s,a) ERs,a	+ y \(\sigma\) \(\sigma\) \((s')\)			
of he have dis,a	) = max Q (5,a)	(by	det, of	Q*(s,a))
	= max (RS, a+ 3) 5'45	$P_{SS}$ , $V(S')$	def. of	Q (S, a)
	(5) < V*(5) YTETT	¥ \$ (by	वर्ध म	V*(s))
By O and Q, Q's	,a)= max(RS,a+) }	FS' (5')		
	< RS,a+ J I Pas		)30,e30	A6 = 622,)
⇒Q*(s,a) ERs,a+y∑		t <sub>f</sub>		
Obv2. Q'(s,a) < Rs,a+ y \( \)				
Yes, by stree 1	Definition of an "Optima	al" Policy (Forma	lly)	
	部例可 <b>以</b> <b>Partial ordering of policies</b> :			
	$\pi \geq \pi'$ if $V$	$\pi(s) \ge V^{\pi'}(s), \forall s$		
• An optimal policy: A policy $\pi^*$ is an optimal policy if it is better than or equal to all other policies, i.e.				
	$\pi^* \geq \pi$ , 1	or all $\pi \in \Pi$		
•	Question: Given an MDP, does s	uch $\pi^*$ always exist?	- [	3
By Obl and Obv	Z, we know Qts	N=RSIA+ YI	o V'(5')	

l

```
1.(b) We need to proof 11T(Q)-T(Q)112=3/11Q-Q1120 VQ,Q'
    11/(Q)-1*(Q)11= max[[]*(Q)](s,a)-[]*(d)](s,a)] (by def. of on-norm)
                         = max (Rs, at y) Psy max Q(s, a))
                                   -(Rs, a+) Z for max Q(s, a)) (by def of T*(Q))
                        = I max | I Pa (max Q(s'a')-max Q(s'a')) |
sia s' s' a' Q(s'a')-max Q(s'a')) |
                        by max operator property we use in lecture 4 PI)
                          \lim_{s \to \infty} \max_{a} \left| \max_{a} \left( R_{s,a} + \gamma \sum_{s'} P_{ss'}^{a} V(s') \right) - \max_{a'} \left( R_{s,a'} + \gamma \sum_{s'} P_{ss'}^{a'} \hat{V}(s') \right) \right| 
 \leq \max_{s} \max_{a} \left| \left( R_{s,a} + \gamma \sum_{s'} P_{ss'}^{a} V(s') \right) - \left( R_{s,a} + \gamma \sum_{s'} P_{ss'}^{a} \hat{V}(s') \right) \right| 
                       = IIR-allo (by def. of on-norm)
2.(a) We need to proof 1/2(V)-1/2(V))=3/1V-V/1 VV, V'
      11 2 (V) - T2 (V') ||= max |[T2V] (S)-[T2V'] (By the def. of so-norm)
                           = max | (Rs+2 (T(-(s))+) Pt, V) (by the def. of

-(Rs+2 (T(-(s))+) Pt, V) | [T2V](S)
                           = I max | PT (V-V') (cancel the common terms)
                           < y max /V-V'/ (: Pss. 13 a transstion matrix)
                           = y 11 V-V 110 (by the def. of 00-norm)
```

216) I obtain Va Stepl. Initialize k=0 and V2,0(S)=0 for all states S Stop Z. Repeating the following until convergence VII, K+1(S) - max (Rs, a+) I PS, Va, (S') +52(1(2(3))) Step3 Return last computed VK I obtain Q's Stepl. Find Vsz Step 7. Return Q (S,a)= RS,a+ 7] I Pa VT (S')

3. 1 Es die Eart (15) [f(s,a)] (Renne natural)  $= \frac{1}{1-y} \sum_{S} (1-y) \sum_{T} y^{T} P(S_{t}=S|S_{0}, T) \sum_{C_{t}} (S_{0}, T) \sum_{C_{t}} (a_{t}|S_{0}) e_{t}(S_{0}) e_{t}(S_{0})$ = 2 1 = 2 f(s,02) P(St=2(20,7) T(a/s) ms) = 2 yt Zf(Stat) P((2, at) E C) M(So) = \(\int\_{\infty} \int\_{\infty} \int\_{\infty

4

# 5. original=mazezd-amaze-vi

## myjoco halfcheetah-rundom-v2'

```
[[.2.8027589e-02 .0.3696302e-02 8.6608730e-02 ...-8.6148446e-03 '....4.0926412e-03 1.5970843e-01 [.....4.0926412e-03 1.5970843e-01 1.7004584e-00 1.7004584e-
```

2D away 科學記表表示