

(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:

$$V^*(s) = \max_a Q^*(s, a) \quad (1)$$

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

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 $Q^{\pi}(s, a) = R_{s,a} + \gamma \sum_{s'} P_{ss'}^a V^{\pi}(s')$)

△ Suppose: ① For any finite MDP, there always exists an π^* , s.t. $\pi^* \geq \pi \forall \pi \in \Pi$.

② A deterministic defined policy $\pi^*(a|s)$ can always be found.

$$\pi^*(a|s) = \begin{cases} 1, & \text{if } a = \arg \max_{a' \in A} Q^*(s, a') \\ 0, & \end{cases}$$

(1)

$$V^*(s) = \max_{\pi \in \Pi} V^{\pi}(s) \quad \text{and already know that } V^{\pi}(s) = \sum_{a \in A} \pi(a|s) \cdot Q^{\pi}(s, a),$$

With a given optimal policy π^* , for each $s \in S$, take action $a \in \arg \max_{a \in A} Q^{\pi^*}(s, a)$.

$$\text{we have } V^*(s) = V^{\pi^*}(s) = \sum_{a \in A} \pi^*(a|s) \cdot Q^{\pi^*}(s, a) \leq \max_{a \in A} Q^*(s, a) \dots (1)$$

Now assume a state s' , an action a' , and a policy ϕ s.t. $V^*(s') < Q^{\phi}(s', a') = \max_{a \in A} Q^*(s', a)$

$$\Rightarrow V^{\phi}(s') = \sum_{a \in A} \phi(a|s') \cdot Q^{\phi}(s, a) = \max_{a \in A} Q^*(s', a)$$

$$V^{\phi}(s') > V^*(s').$$

which contradicts $V^*(s) = \max_{\pi \in \Pi} V^{\pi}(s)$ *

(2)

$$\text{Since we have } Q^{\pi}(s, a) = R_{s,a} + \gamma \sum_{s' \in S} P_{ss'}^a V^{\pi}(s') \quad \text{and} \quad V^*(s) = \max_{a \in A} Q^*(s, a)$$

With a given optimal policy π^* , for each $s \in S$, take action $a \in \arg \max_{a \in A} Q^{\pi^*}(s, a)$.

$$\text{We can derive } Q^{\pi^*}(s, a) = Q^*(s, a) = R_{s,a} + \gamma \sum_{s' \in S} P_{ss'}^a V^{\pi^*}(s') \longrightarrow V^{\pi^*}(s') = V^*(s')$$

$$\Rightarrow Q^*(s, a) = R_{s,a} + \gamma \sum_{s' \in S} P_{ss'}^a \cdot V^*(s')$$

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

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

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

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

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

by hint: $\|T^*(Q) - T^*(Q')\|_\infty = \max_{(s,a)} |[T^*(Q)](s, a) - [T^*(Q')](s, a)|$

(by def.) $= \max_{(s,a)} \left| \left[\cancel{R_{s,a}} + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q(s', a') \right] - \left[\cancel{R_{s,a}} + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q'(s', a') \right] \right|$

$$= \gamma \cdot \max_{(s,a)} \left| \sum_{s'} P_{ss'}^a \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right|$$

$(|a| + |b| > |a+b|) \leq \gamma \cdot \max_{(s,a)} \sum_{s'} P_{ss'}^a \left| \max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right|$

$\leq \gamma \cdot \max_{(s,a)} \sum_{s'} P_{ss'}^a \max_{a'} |Q(s', a') - Q'(s', a')|$

$\leq \gamma \|Q - Q'\|_\infty$

$\max_{x \in X} |f(x) - g(x)| \geq \left| \max_x f(x) - \max_x g(x) \right|$

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 Ω . In this problem, for simplicity, suppose we use the Shannon entropy as our regularizer, i.e., $\Omega(\pi(\cdot|s)) \equiv H(\pi(\cdot|s)) := -\sum_{a \in A} \pi(a|s) \ln \pi(a|s)$. Let us verify a few important properties mentioned in Lecture 4 as follows.

(a) Recall that we introduce the “regularized Bellman expectation operator” T_Ω^π as

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

Please verify that T_Ω^π is a contraction operator in L_∞ norm. (Hint: Try to extend the proof procedure of the contraction property of T^π in Lecture 3)

$$\|T_\Omega^\pi(V) - T_\Omega^\pi(V')\|_\infty = \max_s \left| [T_\Omega^\pi(V)](s) - [T_\Omega^\pi(V')](s) \right|$$

$$= \max_s \left| \left(\cancel{R_s^\pi} + \cancel{\Omega(\pi(\cdot|s))} + \gamma \sum_{s'} P_{ss'}^\pi V(s') \right) - \left(\cancel{R_s^\pi} + \cancel{\Omega(\pi(\cdot|s))} + \gamma \sum_{s'} P_{ss'}^\pi V'(s') \right) \right|$$

$$= \max_s \left| \gamma \sum_{s'} P_{ss'}^\pi (V(s') - V'(s')) \right|$$

$$\leq \gamma \max_s \sum_{s'} P_{ss'}^\pi |V(s') - V'(s')| \leq \gamma \|V - V'\|_\infty$$

(b) Moreover, under regularized MDPs, we study the optimal value functions V_{Ω}^* and optimal Q functions Q_{Ω}^*

and learn the Bellman optimality equations as

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

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

Could you design an iterative algorithm that can obtain V_{Ω}^* and Q_{Ω}^* ? 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)

For $V_{\Omega}^*(s)$:

initialize $V_{\Omega}^*(s) = 0 \quad \forall s \in S$, $\Delta \leftarrow 0$, $v \leftarrow V_{\Omega}^*(s)$

Repeat until $V_{\Omega}^{\pi}(s)$ converge

for $s \in S$:

$V_{\Omega}^*(s) \leftarrow \max_{\pi \in \Pi} R_s^{\pi} + \gamma \mathbb{P}_s^{\pi} V_{\Omega}^*(s')$

For $s \in S$:

$\pi^*(s) \leftarrow \operatorname{argmax}_{\pi \in \Pi} \left(\max_{\pi \in \Pi} R_s^{\pi} + \gamma \mathbb{P}_s^{\pi} \cdot V_{\Omega}^*(s') \right)$

For $Q_{\Omega}^*(s, a)$:

initialize $V_{\Omega}^*(s) = 0 \quad \forall s \in S$, initialize $\pi(s)$ randomly

Repeat until $\pi(s)$ converge:

Repeat until $V_{\Omega}^{\pi}(s)$ converge

for $s \in S$

$V_{\Omega}^*(s) \leftarrow \max_{\pi \in \Pi} R_s^{\pi} + \gamma \mathbb{P}_s^{\pi} V_{\Omega}^*(s')$

For all (s, a) pair

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

$\pi^*(s) = \operatorname{argmax}_{a \in A} (Q_{\Omega}^*(s, a))$

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

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

(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 $d_{\mu}^{\pi_{\theta}}$, which involves a sum of probability over t . Next, try to reorganize the triple summation into the form of the LHS of (8))

$$\mathbb{E}_{\tau \sim P_{\mu}^{\pi_{\theta}}} \left[\sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \sum_{\tau} \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) P_{\mu}^{\pi_{\theta}}(\tau)$$

$$\frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta}}} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [f(s, a)] = \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} d_{\mu}^{\pi_{\theta}}(s) \cdot \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [f(s, a)]$$

$$\because d_{\mu}^{\pi_{\theta}}(s) := \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^{\pi_{\theta}}(s)] = \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^{\pi_{\theta}}(s)] \cdot \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [f(s, a)]$$

$$\because d_{s_0}^{\pi_{\theta}}(s) = \frac{1}{1-\gamma} \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{P}^{\pi_{\theta}}(s_t = s | s_0) = \sum_{s \in \mathcal{S}} \mathbb{E}_{s_0 \sim \mu} \left[\sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{P}^{\pi_{\theta}}(s_t = s | s_0) \right] \cdot \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [f(s, a)]$$

$$\begin{aligned} \text{Triple summation:} &= \sum_{s \in \mathcal{S}} \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{P}_{\mu}^{\pi_{\theta}} \cdot \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) \cdot [f(s_t, a_t)] \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{P}_{\mu}^{\pi_{\theta}} f(s_t, a_t) \\ &= \sum_{\tau} \mathbb{P}_{\mu}^{\pi_{\theta}} \left[\sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \mathbb{E}_{\tau \sim P_{\mu}^{\pi_{\theta}}} \left[\sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] \end{aligned}$$

△ The observation and output of problem 5, are shown in the zip. file "dataset" is a dictionary, observation is ndarray.