

Homework 1–Part I: Planning for MDPs

Submission Guidelines: Please compress all your write-ups (photos/scanned copies are acceptable; please make sure that the electronic files are of good quality and reader-friendly) and source code into one .zip file and submit the compressed file via E3.

Problem 1 (Q-Value Iteration)

(20+15=35 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:

$$V_*(s) = \max_a Q_*(s, a) \quad (1)$$

$$Q_*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V_*(s'). \quad (2)$$

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

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

$$Q_*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \left(\max_{a'} Q_*(s', a') \right) \quad (3)$$

Similar to the value iteration, we can study the *Q-value iteration* by defining the Bellman optimality operator $T^* : \mathbb{R}^{|S||A|} \rightarrow \mathbb{R}^{|S||A|}$ 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 $\|T^*(Q) - T^*(Q')\|_{\infty} = \max_{(s,a)} |[T^*(Q)](s, a) - [T^*(Q')](s, a)|$)

Problem 2 (Distributional Perspective of MDPs)

(20+15=35 points)

Recall that given a policy π , the distributional Bellman operator $B^{\pi} : \mathcal{Z} \rightarrow \mathcal{Z}$ is defined as

$$[B^{\pi}Z](s, a) \stackrel{D}{=} r(s, a) + \gamma P^{\pi}Z(s, a), \quad (5)$$

where $\gamma \in (0, 1)$. In the following subproblems, we would like to show that the B^{π} is a contraction operator in the maximal form of the Wasserstein metric (i.e. \bar{d}_p defined in Lecture 4). For ease of exposition, we further consider the following notations: Given any two random variables U, V with CDFs F_U, F_V , we write $d_p(U, V) := d_p(F_U, F_V)$.

(a) To begin with, show that the Wasserstein metric satisfies the following nice properties: Let U and V be two random variables. Let A be another random variable that is independent of U and V . Let Q be a Bernoulli random variable that is independent of U and V and satisfies $P(Q = 1) = q$:

- (i) $d_p(aU, aV) = |a|d_p(U, V)$, for any $a \in \mathbb{R}$

- (ii) $d_p(A + U, A + V) \leq d_p(U, V)$
- (iii) $d_p(QU, QV) \leq q \cdot d_p(U, V)$

(Hint: For (i), you may first show that $d_p(aU, aV) \leq |a|d_p(U, V)$; For (ii), for any pair of random variables U', V' with $U' \stackrel{D}{=} U$, $V' \stackrel{D}{=} V$, consider some random variable A' that satisfies $A' \stackrel{D}{=} A$ and is independent of U', V' . Then, try to connect $d_p(A' + U', A' + V')$ and $d_p(U, V)$; For (iii), based on each possible joint distribution of U, V , construct one straightforward joint distribution of QU, QV)

(b) By using the result in (a) and the partition lemma (Lemma 1 in [Belleware et al., ICML 2017]), show that B^π is a γ -contraction operator in \bar{d}_p . (Hint: As an intermediate step of the proof, you may need to show that $d_p(B^\pi Z_1(s, a), B^\pi Z_2(s, a)) \leq \gamma \sup_{\bar{s}, \bar{a}} d_p(Z_1(\bar{s}, \bar{a}), Z_2(\bar{s}, \bar{a}))$, for any state-action pair (s, a))

Problem 3 (Implementing Policy Iteration and Value Iteration)

(40 points)

In this problem, we will implement policy iteration and value 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/>. Read through `policy_and_value_iteration.py` and then implement the two functions `policy_iteration` and `value_iteration` (Note: please set $\gamma = 0.9$ and the termination criterion $\varepsilon = 10^{-3}$. Moreover, you could use either Taxi-v2 or Taxi-v3 environment. Note that discrepancy = 0 is a necessary condition 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).