

CS234: Reinforcement Learning – Problem Session #3

Winter 2021-2022

Problem 1

Suppose we have an infinite-horizon, discounted MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma \rangle$ with a finite state-action space, $|\mathcal{S} \times \mathcal{A}| < \infty$ and $0 \leq \gamma < 1$. For any two arbitrary sets \mathcal{X} and \mathcal{Y} , we denote the class of all functions mapping from \mathcal{X} to \mathcal{Y} as $\{\mathcal{X} \rightarrow \mathcal{Y}\} \triangleq \{f \mid f : \mathcal{X} \rightarrow \mathcal{Y}\}$. In the questions that follow, let $Q, Q' \in \{\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\}$ be any two arbitrary action-value functions and consider any fixed state $s \in \mathcal{S}$. Without loss of generality, you may assume that $Q(s, a) \geq Q'(s, a), \forall (s, a) \in \mathcal{S} \times \mathcal{A}$.

1. Prove that $\left| \max_{a \in \mathcal{A}} Q(s, a) - \max_{a' \in \mathcal{A}} Q'(s, a') \right| \leq \max_{a \in \mathcal{A}} |Q(s, a) - Q'(s, a)|$.

2. Prove that $\left| \min_{a \in \mathcal{A}} Q(s, a) - \min_{a' \in \mathcal{A}} Q'(s, a') \right| \leq \max_{a \in \mathcal{A}} |Q(s, a) - Q'(s, a)|$.

3. Prove that $\left| \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} Q(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a' \in \mathcal{A}} Q'(s, a') \right| \leq \max_{a \in \mathcal{A}} |Q(s, a) - Q'(s, a)|$.

4. Prove that, for any parameter $\omega \in \mathbb{R}$,¹

$$\left| \frac{1}{\omega} \log \left(\frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \exp(\omega \cdot Q(s, a)) \right) - \frac{1}{\omega} \log \left(\frac{1}{|\mathcal{A}|} \sum_{a' \in \mathcal{A}} \exp(\omega \cdot Q'(s, a')) \right) \right| \leq \max_{a \in \mathcal{A}} |Q(s, a) - Q'(s, a)|.$$

Hint: define and introduce a $\Delta(a) = Q(s, a) - Q'(s, a)$ term where $a \in \mathcal{A}$.

¹For any $x \in \mathbb{R}$, $\exp(x) = e^x$ and all logarithms are base e .

The remainder of this question focuses on Algorithm 1, which takes as input an operator

$$\otimes : \{\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\} \rightarrow \{\mathcal{S} \rightarrow \mathbb{R}\}$$

that adheres to the following property²:

$$\|\otimes Q - \otimes Q'\|_\infty \leq \|Q - Q'\|_\infty, \quad \forall Q, Q' \in \{\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\}. \quad (1)$$

Algorithm 1:

Data: Finite MDP \mathcal{M} , Operator \otimes satisfying Equation 1
Initialize $V_0(s) = 0, \forall s \in \mathcal{S}$ ▷ Initial value function estimate
Initialize $k = 1$ ▷ Iteration counter
while *not converged* **do**
 for *each state* $s \in \mathcal{S}$ **do**
 $V_k(s) = \otimes_{a \in \mathcal{A}} \left(\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{T}(s' | s, a) V_{k-1}(s') \right)$.
 end
 $k = k + 1$
end
Return V_k

5. For any value function $V \in \{\mathcal{S} \rightarrow \mathbb{R}\}$, define the operator $\mathcal{B} : \{\mathcal{S} \rightarrow \mathbb{R}\} \rightarrow \{\mathcal{S} \rightarrow \mathbb{R}\}$ as follows:

$$\mathcal{B}V(s) = \otimes_{a \in \mathcal{A}} \left(\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{T}(s' | s, a) V(s') \right),$$

where \otimes satisfies Equation 1. Prove that \mathcal{B} is a γ -contraction with respect to the L_∞ -norm.

²As always, $\|\cdot\|_\infty$ denotes the L_∞ -norm.

6. Let $\otimes_1, \otimes_2 : \{\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\} \rightarrow \{\mathcal{S} \rightarrow \mathbb{R}\}$ be two operators satisfying Equation 1. Prove that, for any $0 \leq \lambda \leq 1$,

$$\otimes_\lambda = \lambda \otimes_1 + (1 - \lambda) \otimes_2$$

also satisfies Equation 1.

7. For any $0 \leq \varepsilon \leq 1$, define your own operator $\otimes_\varepsilon : \{\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\} \rightarrow \{\mathcal{S} \rightarrow \mathbb{R}\}$ and prove that running Algorithm 1 with your \otimes_ε returns the value function associated with the ε -greedy optimal policy (where the optimal policy maximizes the expected sum of future discounted rewards).

Problem 2

Consider an infinite-horizon, discounted MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma \rangle$ where $\gamma \in [0, 1)$ and the state-action space is finite ($|\mathcal{S} \times \mathcal{A}| < \infty$). For any policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, recall that the discounted stationary state distribution is defined for any state $s \in \mathcal{S}$ as

$$d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \beta_t^\pi(s_t = s),$$

where $\beta_t^\pi(s_t = s)$ denotes the probability that the (random) state s_t encountered by policy π at timestep t is equal to s . Let $\beta \in \Delta(\mathcal{S})$ be an initial state distribution such that $\beta_0^\pi = \beta(s)$ for all policies π and any state $s \in \mathcal{S}$. Prove that for any state $s' \in \mathcal{S}$,

$$d^\pi(s') = (1 - \gamma)\beta(s') + \gamma \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \mathcal{T}(s' | s, a) \pi(a | s) d^\pi(s).$$