# CS234: Reinforcement Learning – Problem Session #2

## Winter 2022-2023

## Problem 1

Consider an infinite-horizon, discounted MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma \rangle$.

1. Define the maximal reward $R_{\mathrm{MAX}} = \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathcal{R}(s,a)$ and show that, for any policy $\pi : \mathcal{S} \to \mathcal{A}$,

$$V^\pi(s) \le \frac{R_{\mathrm{MAX}}}{1 - \gamma}, \qquad \forall s \in \mathcal{S}.$$

   Solution:

$$\begin{aligned}
V^\pi(s) &= \mathbb{E}\left[\sum_{t=0}^\infty \gamma^t \mathcal{R}(s_t, a_t) \mid s_0 = s, \pi\right] \\
&\le \mathbb{E}\left[\sum_{t=0}^\infty \gamma^t R_{\mathrm{MAX}} \mid s_0 = s, \pi\right] \\
&= R_{\mathrm{MAX}} \cdot \mathbb{E}\left[\sum_{t=0}^\infty \gamma^t \mid s_0 = s, \pi\right] \\
&= \frac{R_{\mathrm{MAX}}}{1 - \gamma}.
\end{aligned}$$

2. Consider a second MDP $\widehat{\mathcal{M}} = \langle \mathcal{S}, \mathcal{A}, \widehat{\mathcal{R}}, \widehat{\mathcal{T}}, \gamma \rangle$ and define the constant $V_{\mathrm{MAX}} = \frac{R_{\mathrm{MAX}}}{1-\gamma}$. We will use subscripts to distinguish between arbitrary value functions $V_\mathcal{M}$ and $V_{\widehat{\mathcal{M}}}$ of MDPs $\mathcal{M}$ and $\widehat{\mathcal{M}}$, respectively. Suppose we have two constants $\varepsilon_1, \varepsilon_2 > 0$ such that

$$\max_{s,a \in \mathcal{S} \times \mathcal{A}} |\mathcal{R}(s,a) - \widehat{\mathcal{R}}(s,a)| \le \varepsilon_1 \qquad \max_{s,a \in \mathcal{S} \times \mathcal{A}} \sum_{s' \in \mathcal{S}} |\mathcal{T}(s'|s,a) - \widehat{\mathcal{T}}(s'|s,a)| \le \varepsilon_2.$$

   For any policy $\pi : \mathcal{S} \to \mathcal{A}$, show that

$$||V_\mathcal{M}^\pi - V_{\widehat{\mathcal{M}}}^\pi||_\infty \le \frac{\varepsilon_1 + \gamma \varepsilon_2 V_{\mathrm{MAX}}}{(1 - \gamma)}.$$

Solution:

$$\|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} = \max_{s \in \mathcal{S}} \left| V_{\mathcal{M}}^{\pi}(s) - V_{\widehat{\mathcal{M}}}^{\pi}(s) \right|$$

$$= \max_{s \in \mathcal{S}} \left| \mathcal{R}(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) V_{\mathcal{M}}^{\pi}(s') - \widehat{\mathcal{R}}(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} \widehat{\mathcal{T}}(s'|s, \pi(s)) V_{\widehat{\mathcal{M}}}^{\pi}(s') \right|$$

$$\leq \max_{s \in \mathcal{S}} \left| \mathcal{R}(s, \pi(s)) - \widehat{\mathcal{R}}(s, \pi(s)) \right| + \left| \gamma \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) V_{\mathcal{M}}^{\pi}(s') - \gamma \sum_{s' \in \mathcal{S}} \widehat{\mathcal{T}}(s'|s, \pi(s)) V_{\widehat{\mathcal{M}}}^{\pi}(s') \right|$$

$$\leq \varepsilon_1 + \max_{s \in \mathcal{S}} \gamma \left| \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) V_{\mathcal{M}}^{\pi}(s') - \sum_{s' \in \mathcal{S}} \widehat{\mathcal{T}}(s'|s, \pi(s)) V_{\widehat{\mathcal{M}}}^{\pi}(s') \right|$$

$$= \varepsilon_1 + \max_{s \in \mathcal{S}} \gamma \left| \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) V_{\mathcal{M}}^{\pi}(s') - \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) V_{\widehat{\mathcal{M}}}^{\pi}(s') + \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) V_{\widehat{\mathcal{M}}}^{\pi}(s') - \sum_{s' \in \mathcal{S}} \widehat{\mathcal{T}}(s'| \right.$$

$$\leq \varepsilon_1 + \max_{s \in \mathcal{S}} \gamma \left| \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) \left( V_{\mathcal{M}}^{\pi}(s') - V_{\widehat{\mathcal{M}}}^{\pi}(s') \right) \right| + \gamma \left| \sum_{s' \in \mathcal{S}} V_{\widehat{\mathcal{M}}}^{\pi}(s') \left( \mathcal{T}(s'|s, \pi(s)) - \widehat{\mathcal{T}}(s'|s, \pi(s)) \right) \right|$$

$$\leq \varepsilon_1 + \max_{s \in \mathcal{S}} \gamma \left| \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} \right| + \gamma \left| \sum_{s' \in \mathcal{S}} V_{\widehat{\mathcal{M}}}^{\pi}(s') \left( \mathcal{T}(s'|s, \pi(s)) - \widehat{\mathcal{T}}(s'|s, \pi(s)) \right) \right|$$

$$= \varepsilon_1 + \max_{s \in \mathcal{S}} \gamma \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} \left| \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, \pi(s)) \right| + \gamma \left| \sum_{s' \in \mathcal{S}} V_{\widehat{\mathcal{M}}}^{\pi}(s') \left( \mathcal{T}(s'|s, \pi(s)) - \widehat{\mathcal{T}}(s'|s, \pi(s)) \right) \right|$$

$$= \varepsilon_1 + \gamma \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} \left| 1 \right| + \max_{s \in \mathcal{S}} \gamma \left| \sum_{s' \in \mathcal{S}} V_{\widehat{\mathcal{M}}}^{\pi}(s') \left( \mathcal{T}(s'|s, \pi(s)) - \widehat{\mathcal{T}}(s'|s, \pi(s)) \right) \right|$$

$$= \varepsilon_1 + \gamma \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} + \max_{s \in \mathcal{S}} \gamma \left| \sum_{s' \in \mathcal{S}} V_{\widehat{\mathcal{M}}}^{\pi}(s') \left( \mathcal{T}(s'|s, \pi(s)) - \widehat{\mathcal{T}}(s'|s, \pi(s)) \right) \right|$$

$$\leq \varepsilon_1 + \max_{s \in \mathcal{S}} \gamma \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} + \max_{s \in \mathcal{S}} \gamma \left| \sum_{s' \in \mathcal{S}} (V_{\text{MAX}}) \left( \mathcal{T}(s'|s, \pi(s)) - \widehat{\mathcal{T}}(s'|s, \pi(s)) \right) \right|$$

$$= \varepsilon_1 + \gamma \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} + \max_{s \in \mathcal{S}} \gamma V_{\text{MAX}} \left| \sum_{s' \in \mathcal{S}} \left( \mathcal{T}(s'|s, \pi(s)) - \widehat{\mathcal{T}}(s'|s, \pi(s)) \right) \right|$$

$$\leq \varepsilon_1 + \gamma \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} + \gamma \varepsilon_2 V_{\text{MAX}}$$

$$\implies \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} - \gamma \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} \leq \varepsilon_1 + \gamma \varepsilon_2 V_{\text{MAX}}$$

$$(1 - \gamma) \|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} \leq \varepsilon_1 + \gamma \varepsilon_2 V_{\text{MAX}}$$

$$\|V_{\mathcal{M}}^{\pi} - V_{\widehat{\mathcal{M}}}^{\pi}\|_{\infty} \leq \frac{\varepsilon_1 + \gamma \varepsilon_2 V_{\text{MAX}}}{(1 - \gamma)}$$

The cited result is known as the simulation lemma [Kearns and Singh, 2002]. The simulation lemma tells us that when we can recover an accurate approximation to the true reward function and transition function of our MDP, the return obtained by policies in the approximate MDP will be close to those of the original MDP.

# References

Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2-3):209–232, 2002.