

# CS234: Reinforcement Learning – Problem Session #1

Winter 2022-2023

## Problem 1

Consider an infinite-horizon, discounted MDP  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma \rangle$ . As usual, for any policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ , the value function induced by  $\pi$  is defined as

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t) \mid s_0 = s, \pi \right].$$

1. For an arbitrary  $Z \in \mathbb{N}$ , consider learning with  $Z + 1$  distinct discount factors  $\gamma_0, \gamma_1, \dots, \gamma_Z$  where the final discount factor matches that of the MDP  $\mathcal{M}$ ,  $\gamma_Z = \gamma$ . Letting  $[Z] \triangleq \{1, 2, \dots, Z\}$  denote the index set, we define the following functions for any policy  $\pi$ :

$$V_{\gamma_z}^\pi = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma_z^t \mathcal{R}(s_t, a_t) \mid s_0 = s, \pi \right] \quad W_z^\pi = V_{\gamma_z}^\pi - V_{\gamma_{z-1}}^\pi, \quad \forall z \in [Z]$$

where  $W_0 = V_{\gamma_0}^\pi$ .

**Solution:** The results of this part were derived by [Romoff et al. \[2019\]](#) who both empirically and theoretically study the benefits of decomposing a single monolithic value function across multiple time-scales through smaller discount factors.

- (a) For any  $z \in [Z]$ ; any policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ ; and any  $s \in \mathcal{S}$ , write an expression for  $V_{\gamma_z}^\pi(s)$  exclusively in terms of  $\{W_0^\pi, W_1^\pi, \dots, W_Z^\pi\}$ .

**Solution:** From the relationships defined above, we can see that

$$V_{\gamma_z}^\pi(s) = \sum_{i=0}^z W_i^\pi(s).$$

- (b) Show that  $W_z^\pi$  obeys the following Bellman equation for any  $z \in [Z]$  and  $s \in \mathcal{S}$ :

$$W_z^\pi(s) = \mathbb{E}_{\substack{a \sim \pi(\cdot|s) \\ s' \sim \mathcal{T}(\cdot|s,a)}} \left[ (\gamma_z - \gamma_{z-1}) V_{\gamma_{z-1}}^\pi(s') + \gamma_z W_z^\pi(s') \right]$$

**Solution:** Just by expanding the corresponding Bellman equations for  $V_{\gamma_z}^\pi$  and  $V_{\gamma_{z-1}}^\pi$ , we have

$$\begin{aligned} W_z^\pi(s) &= V_{\gamma_z}^\pi - V_{\gamma_{z-1}}^\pi \\ &= \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \mathcal{R}(s, a) + \gamma_z \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s,a)} [V_{\gamma_z}^\pi(s')] - \mathcal{R}(s, a) - \gamma_{z-1} \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s,a)} [V_{\gamma_{z-1}}^\pi(s')] \right] \\ &= \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \gamma_z \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s,a)} [V_{\gamma_z}^\pi(s')] - \gamma_{z-1} \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s,a)} [V_{\gamma_{z-1}}^\pi(s')] \right] \\ &= \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \gamma_z \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s,a)} [W_z^\pi(s') + V_{\gamma_{z-1}}^\pi(s')] - \gamma_{z-1} \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s,a)} [V_{\gamma_{z-1}}^\pi(s')] \right] \\ &= \mathbb{E}_{\substack{a \sim \pi(\cdot|s) \\ s' \sim \mathcal{T}(\cdot|s,a)}} \left[ (\gamma_z - \gamma_{z-1}) V_{\gamma_{z-1}}^\pi(s') + \gamma_z W_z^\pi(s') \right]. \end{aligned}$$

2. Let  $\gamma, \beta \in [0, 1)$  be two discount factors such that  $\beta \leq \gamma$ . Let  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$  be an arbitrary policy that induces value functions  $V_\gamma^\pi$  and  $V_\beta^\pi$  under the two discount factors, respectively. Similarly, define the Bellman operators

$$\begin{aligned}\mathcal{B}_\gamma^\pi V(s) &= \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathcal{R}(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s, a)} [V(s')]] \\ \mathcal{B}_\beta^\pi V(s) &= \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathcal{R}(s, a) + \beta \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s, a)} [V(s')]] .\end{aligned}$$

With the reward upper bound  $R_{\text{MAX}} = \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \mathcal{R}(s, a)$ , prove that

$$\|V_\gamma^\pi - V_\beta^\pi\|_\infty \leq \frac{(\gamma - \beta)R_{\text{MAX}}}{(1 - \gamma)(1 - \beta)}.$$

**Solution:** This result is given as Theorem 2 of [Petric and Scherrer, 2008] and highlights the approximation error that can occur by using a smaller discount factor  $\beta$  than that of the true MDP,  $\gamma$ .

$$\begin{aligned}\|V_\gamma^\pi - V_\beta^\pi\|_\infty &= \|\mathcal{B}_\gamma^\pi V_\gamma^\pi - \mathcal{B}_\beta^\pi V_\beta^\pi\|_\infty \\ &= \|\mathcal{B}_\gamma^\pi V_\gamma^\pi - \mathcal{B}_\beta^\pi V_\gamma^\pi + \mathcal{B}_\beta^\pi V_\gamma^\pi - \mathcal{B}_\beta^\pi V_\beta^\pi\|_\infty \\ &\leq \|\mathcal{B}_\gamma^\pi V_\gamma^\pi - \mathcal{B}_\beta^\pi V_\gamma^\pi\|_\infty + \|\mathcal{B}_\beta^\pi V_\gamma^\pi - \mathcal{B}_\beta^\pi V_\beta^\pi\|_\infty \\ &\leq \|\mathcal{B}_\gamma^\pi V_\gamma^\pi - \mathcal{B}_\beta^\pi V_\gamma^\pi\|_\infty + \beta \|V_\gamma^\pi - V_\beta^\pi\|_\infty \\ &= \max_{s \in \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)} [\mathcal{R}(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s, a)} [V_\gamma^\pi(s')] - \mathcal{R}(s, a) - \beta \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s, a)} [V_\gamma^\pi(s')]]| + \beta \|V_\gamma^\pi - V_\beta^\pi\|_\infty \\ &= \max_{s \in \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)} [\gamma \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s, a)} [V_\gamma^\pi(s')] - \beta \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s, a)} [V_\gamma^\pi(s')]]| + \beta \|V_\gamma^\pi - V_\beta^\pi\|_\infty \\ &= \max_{s \in \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)} [(\gamma - \beta) \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s, a)} [V_\gamma^\pi(s')]]| + \beta \|V_\gamma^\pi - V_\beta^\pi\|_\infty \\ &\leq \max_{s \in \mathcal{S}} |\mathbb{E}_{a \sim \pi(\cdot|s)} \left[ (\gamma - \beta) \mathbb{E}_{s' \sim \mathcal{T}(\cdot|s, a)} \left[ \frac{R_{\text{MAX}}}{(1 - \gamma)} \right] \right]| + \beta \|V_\gamma^\pi - V_\beta^\pi\|_\infty \\ &= \frac{(\gamma - \beta)R_{\text{MAX}}}{(1 - \gamma)} + \beta \|V_\gamma^\pi - V_\beta^\pi\|_\infty \\ \implies (1 - \beta) \|V_\gamma^\pi - V_\beta^\pi\|_\infty &\leq \frac{(\gamma - \beta)R_{\text{MAX}}}{(1 - \gamma)} \\ \|V_\gamma^\pi - V_\beta^\pi\|_\infty &\leq \frac{(\gamma - \beta)R_{\text{MAX}}}{(1 - \gamma)(1 - \beta)}\end{aligned}$$

3. Let  $\alpha, \gamma \in [0, 1)$  be two discount factors such that  $\gamma \leq \alpha$ . Consider a new MDP  $\mathcal{M}' = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}', \mathcal{R}, \alpha \rangle$  with a different transition function  $\mathcal{T}' : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  defined for  $\lambda \in [0, 1]$  as

$$\mathcal{T}'(s' | s, a) = (1 - \lambda)\mathcal{T}(s' | s, a) + \lambda \mathbb{1}(s = s'), \quad \forall (s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}.$$

In words, the new transition function  $\mathcal{T}'$  follows the transitions of the original MDP  $\mathcal{T}$  with probability  $(1 - \lambda)$  and takes a self-looping transition with probability  $\lambda$ . We will use subscripts to distinguish between value functions of  $\mathcal{M}$  versus those of  $\mathcal{M}'$ .

Assuming that both  $\mathcal{M}$  and  $\mathcal{M}'$  are tabular, recall the matrix form of the Bellman equations for any policy  $\pi$ :

$$V_{\mathcal{M}}^\pi = (I - \gamma \mathcal{T}^\pi)^{-1} \mathcal{R}^\pi \quad V_{\mathcal{M}'}^\pi = (I - \alpha \mathcal{T}'^\pi)^{-1} \mathcal{R}^\pi,$$

where

$$\mathcal{R}^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathcal{R}(s, a)] \quad \mathcal{T}^\pi(s' | s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathcal{T}(s' | s, a)] \quad \mathcal{T}'^\pi(s' | s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathcal{T}'(s' | s, a)]$$

Solution: The results of this question are proven as part of Theorem 1 in [Jiang et al., 2015].

(a) Give a value of  $\lambda$  such that, for any policy  $\pi$ ,

$$V_{\mathcal{M}'}^\pi = \frac{1-\gamma}{1-\alpha} \cdot V_{\mathcal{M}}^\pi.$$

Solution: We can write the transition matrix in the new MDP  $\mathcal{M}'$  induced by any policy  $\pi$  as

$$\mathcal{T}'^\pi = (1-\lambda)\mathcal{T}^\pi + \lambda I,$$

where  $I$  is the  $|\mathcal{S}| \times |\mathcal{S}|$  identity matrix. So, substituting in directly, we have

$$\begin{aligned} V_{\mathcal{M}'}^\pi &= (I - \alpha \mathcal{T}'^\pi)^{-1} \mathcal{R}^\pi \\ &= (I - \alpha((1-\lambda)\mathcal{T}^\pi + \lambda I))^{-1} \mathcal{R}^\pi \\ &= ((1-\alpha\lambda)I - \alpha(1-\lambda)\mathcal{T}^\pi)^{-1} \mathcal{R}^\pi \\ &= \left( (1-\alpha\lambda) \left( I - \frac{\alpha(1-\lambda)}{1-\alpha\lambda} \mathcal{T}^\pi \right) \right)^{-1} \mathcal{R}^\pi \\ &= \frac{1}{1-\alpha\lambda} \left( I - \frac{\alpha(1-\lambda)}{1-\alpha\lambda} \mathcal{T}^\pi \right)^{-1} \mathcal{R}^\pi. \end{aligned}$$

We can compute the required value of  $\lambda$  as

$$\frac{\alpha(1-\lambda)}{1-\alpha\lambda} = \gamma \implies \lambda = \frac{\alpha-\gamma}{\alpha(1-\gamma)},$$

which means

$$\frac{1}{1-\alpha\lambda} = \frac{1}{1 - \frac{\alpha-\gamma}{(1-\gamma)}} = \frac{1-\gamma}{1-\gamma-\alpha+\gamma} = \frac{1-\gamma}{1-\alpha}.$$

Substituting back in to the earlier equation yields

$$\begin{aligned} V_{\mathcal{M}'}^\pi &= \frac{1}{1-\alpha\lambda} \left( I - \frac{\alpha(1-\lambda)}{1-\alpha\lambda} \mathcal{T}^\pi \right)^{-1} \mathcal{R}^\pi \\ &= \frac{1-\gamma}{1-\alpha} (I - \gamma \mathcal{T}^\pi)^{-1} \mathcal{R}^\pi \\ &= \frac{1-\gamma}{1-\alpha} \cdot V_{\mathcal{M}}^\pi. \end{aligned}$$

(b) If  $\pi^*$  is the optimal policy of MDP  $\mathcal{M}$ , prove that  $\pi^*$  is also optimal in  $\mathcal{M}'$ .

Solution: By definition of the optimal policy, we know that  $\pi^*$  obeys the following inequality for any other policy  $\pi$ :

$$V_{\mathcal{M}}^{\pi^*}(s) \geq V_{\mathcal{M}}^\pi(s), \quad \forall s \in \mathcal{S}.$$

Since  $\frac{1-\gamma}{1-\alpha} > 0$ , we can scale both sides to get

$$\frac{1-\gamma}{1-\alpha} \cdot V_{\mathcal{M}}^{\pi^*}(s) \geq \frac{1-\gamma}{1-\alpha} \cdot V_{\mathcal{M}}^\pi(s), \quad \forall s \in \mathcal{S}.$$

Applying this previous part, we see that for any other policy  $\pi$ ,

$$V_{\mathcal{M}'}^{\pi^*}(s) \geq V_{\mathcal{M}'}^{\pi}(s), \quad \forall s \in \mathcal{S}.$$

Thus, by definition,  $\pi^*$  is also the optimal policy in MDP  $\mathcal{M}'$ . This result illustrates that, for any MDP with a particular discount factor, there exists a transition function for another MDP with a larger discount factor such that the two MDPs have the same optimal policy.

## References

- Nan Jiang, Alex Kulesza, Satinder Singh, and Richard Lewis. The dependence of effective planning horizon on model accuracy. In *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pages 1181–1189. Citeseer, 2015.
- Marek Petrik and Bruno Scherrer. Biasing approximate dynamic programming with a lower discount factor. *Advances in Neural Information Processing Systems*, 21, 2008.
- Joshua Romoff, Peter Henderson, Ahmed Touati, Emma Brunskill, Joelle Pineau, and Yann Ollivier. Separating value functions across time-scales. In *International Conference on Machine Learning*, pages 5468–5477. PMLR, 2019.