# Mathematics of Reinforcement Learning

## Exercise Class 8

**Exercise 1.** Let $\gamma < 1$, $\pi$ an arbitrary policy and define an operator $T^\pi$ acting on functions $w\colon S \to \mathbb{R}$ by

$$T^\pi[w](s) := \sum_{s' \in S,\, a \in A} \left[ r(s,a) + \gamma w(s') \right] \pi(a \mid s) p(s' \mid s,a), \quad \text{for } s \in S.$$

Show that $T^\pi$ has a unique fixed point $W\colon S \to \mathbb{R}$, that is, $W(s) = T^\pi[W](s)$ for all $s \in S$.

**Exercise 2.** Let $M = (S, A, D, p, r, \gamma)$ be a Markov Decision Model, with $\gamma < 1$. Show that there exists a Markov Decision Model $\tilde{M} = (S, A, D, \tilde{p}, \tilde{r}, \gamma)$ such that

$$\sup_{\pi \in \Pi^M \ \varepsilon\text{-soft}} V^\pi(s) = \sup_{\pi \in \Pi^{\tilde{M}}} V^\pi(s), \quad \text{for all } s \in S, \qquad \text{\color{red}{Also works for ε-greedy?}}$$

where $\Pi^M$ and $\Pi^{\tilde{M}}$ denote the set of all policies in the Markov Decision Models $M$ and $\tilde{M}$.

*Hints:* Start by defining a new transition probability function $\tilde{p}$ that incorporates the $\varepsilon$-softness into the new MDM and adjust the reward function. Show that there exists a transformation of the policies such that the value function remains invariant with respect to these changes. For the last point, use the results of exercise 1. ◇

**Definition 1** ($\varepsilon$-soft optimal)**.** An $\varepsilon$-soft policy $\pi^*$ is called *$\varepsilon$-soft optimal* if

$$V^{\pi^*}(s) = \sup_{\pi \ \varepsilon\text{-soft}} V^\pi(s) =: \tilde{V}^*(s), \quad \text{for all } s \in S.$$

**Exercise 3.** Let $\gamma < 1$ and $\pi_0$ be an arbitrary $\varepsilon$-soft policy and $\{\pi_n\}_{n \in \mathbb{N}} \subseteq \Pi$ be a sequence of $\varepsilon$-soft policies, where $\pi_n$ is chosen to be $\varepsilon$-greedy with respect to $Q^{\pi_{n-1}}$, for all $n > 0$. Show that for some $N \in \mathbb{N}$, for all $m \geq N$, the policy $\pi_m$ is $\varepsilon$-soft optimal.
*Hint:* Use exercise 2. ◇