

Reinforcement learning - HW1

Notations: A MDP (S, A, p, γ)

$$\forall V: S \rightarrow \mathbb{R}, \quad T^* V(s) = \max_a \left[r(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V(s') \right] \quad \text{(Bellman operator optimal)}$$

$$\forall V: S \rightarrow \mathbb{R}, \quad T^\pi V(s) = r(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' | s, \pi(s)) V(s') \quad \text{(Bellman operator for a policy } \pi)$$

Exercise 1: 1) V^* is the optimal value function, hence $\forall s \in S, T^* V^*(s) = V^*(s)$.

So V^* belongs to the feasible set of the optimization problem and a solution exists.

Let us take \tilde{V} a solution for the optimization problem. We have $\arg \min_{V \geq T^* V} \sum_{s \in S} V(s) = \tilde{V}$

By hypothesis, we have that $\sum_{s \in S} T^* \tilde{V}(s) \leq \sum_{s \in S} \tilde{V}(s)$ but since \tilde{V} is the solution of the problem, we have $T^* \tilde{V} = \tilde{V}$. Hence \tilde{V} is fixed point of the Bellman operator and $\tilde{V} = V^*$ since the fixed point is unique.

$V^* \text{ is the unique solution of the problem } \min_{V \geq T^* V} \sum_{s \in S} V(s)$

2) Let $T_S = (p(s' | s, a))_{(s, a) \in S \times A}$ the matrix representation of the probability kernel transition.
et $V = (V(s))_{s \in S}$

Let us pose as well, $\mathbb{1} = (\mathbf{1})_{|S|}$ et $\bar{r} = (r(s, a))_{a \in A}$, then it comes the

following problem equivalence:

$$\min_{V \geq T^* V} \sum_{s \in S} V(s) \quad \Leftrightarrow$$

$$\begin{cases} \min_{\mathbb{R}^{|S|}} \mathbb{1}^T V \\ (\text{Id} - \gamma T_S) V - \bar{r} \geq 0 \end{cases}$$

Exercise 2 - $V: S \mapsto \mathbb{R}$ a function. $\pi: S \rightarrow \mathcal{A}$ the greedy policy w.r.t the value function V :

$$\pi(s) \in \operatorname{argmax}_a \left[r(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V(s') \right]$$

1.a) We obviously have $T^*V \geq T^\pi V$.

Reciprocally, $\forall s, a \in S \times \mathcal{A}$, $T^\pi V(s) \geq r(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V(s')$

hence $\forall s$, $T^\pi V(s) \geq \max_a \left[r(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V(s') \right] = T^*V(s)$

As a conclusion, $\boxed{T^*V = T^\pi V}$.

1.b) $\|V^* - V^\pi\|_\infty = \|V^* - T^\pi \hat{V} + T^\pi \hat{V} - V^\pi\|_\infty$ but since V^* is a fixed point of T^* and by Minkowski inequality

$$\|V^* - V^\pi\|_\infty \leq \|T^*V^* - T^*\hat{V}\|_\infty + \|T^\pi \hat{V} - T^\pi V^\pi\|_\infty \quad \text{and } T^*, T^\pi \text{ are } \gamma\text{-contraction}$$

$$\|V^* - V^\pi\|_\infty \leq \gamma [\|V^* - \hat{V}\|_\infty + \|\hat{V} - V^\pi\|_\infty]$$

$$\|V^* - V^\pi\|_\infty \leq \gamma [\|V^* - \hat{V}\|_\infty + \|\hat{V} - V^*\|_\infty + \|V^* - V^\pi\|_\infty]$$

hence we have the result;

$$\boxed{\|V^* - V^\pi\|_\infty \leq \frac{\gamma}{1-\gamma} \|V^* - \hat{V}\|_\infty}$$

1.c) i) If $V^* = \hat{V}$ then $\|V^* - \hat{V}\|_\infty = 0$ and since 1.b we have

$$\|V^* - V^\pi\|_\infty = 0 \quad \text{so } \forall s \in S, V^*(s) = V^\pi(s) \Rightarrow V^* = V^\pi \Rightarrow \boxed{\hat{V} = V^\pi}$$

ii) Reciprocally, $\hat{V} = V^\pi$ hence $V^\pi = T^*V^\pi$, i.e. V^π is a fixed point of T^* . by definition, the fixed point is unique so $\boxed{V^\pi = V^*}$.

Conclusion: $\boxed{V^\pi = V^* \iff V^* = \hat{V}}$

2-a) - $Q: S \times A \rightarrow \mathbb{R}$ a function and the greedy policy $\tilde{\pi}: S \rightarrow A$

$$\tilde{\pi}(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$$

- $Q_{\tilde{\pi}}: S \times A \rightarrow \mathbb{R}$

$$Q_{\tilde{\pi}}(s, a) = r(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V_{\tilde{\pi}}(s')$$

- Q^* the optimal value function and $\epsilon = \sup_s [Q^*(s, \pi^*(s)) - Q_{\tilde{\pi}}(s, \tilde{\pi}(s))]$

Since π^* is the optimal policy, $V^* \geq V_{\tilde{\pi}} \Rightarrow \forall s \in S, Q^*(s, \pi^*(s)) \geq Q_{\tilde{\pi}}(s, \tilde{\pi}(s))$

hence $\boxed{\epsilon \geq 0}$

2-b) Let us pose $\forall (s, a) \in S \times A$, $Q^*(s, a) - \|Q^* - Q\|_{\infty} \leq Q(s, a) \leq Q^*(s, a) + \|Q^* - Q\|_{\infty}$

Since $\tilde{\pi}$ is the greedy policy w.r.t Q , we have $\forall s \in S, Q(s, \pi^*(s)) \leq Q(s, \tilde{\pi}(s))$

combining the inequalities: $r(s, \pi^*(s)) + \gamma \sum_{s' \in S} p(s' | s, \pi^*(s)) (Q(s', \pi^*(s')) - \epsilon) \leq r(s, \tilde{\pi}(s)) + \gamma \sum_{s' \in S} p(s' | s, \tilde{\pi}(s)) (Q(s', \tilde{\pi}(s')) - \epsilon)$

i.e. ; $Q^*(s, \pi^*(s)) - \|Q^* - Q\|_{\infty} \leq Q^*(s, \tilde{\pi}(s)) + \|Q^* - Q\|_{\infty}$

Then, $Q^*(s, \pi^*(s)) - Q_{\tilde{\pi}}(s, \tilde{\pi}(s)) \leq 2\|Q^* - Q\|_{\infty} + Q^*(s, \tilde{\pi}(s)) - Q_{\tilde{\pi}}(s, \tilde{\pi}(s))$

by definition of $\epsilon \in \operatorname{argmax}_s Q^*(s, \tilde{\pi}(s)) - Q_{\tilde{\pi}}(s, \tilde{\pi}(s))$ $\epsilon \leq 2\|Q^* - Q\|_{\infty} + \gamma \sum_{s' \in S} p(s' | s, \tilde{\pi}(s)) [V^*(s') - V^{\tilde{\pi}}(s')]$

$$\epsilon \leq 2\|Q^* - Q\|_{\infty} + \gamma \sum_{s' \in S} p(s' | s, \tilde{\pi}(s)) [Q^*(s', \pi^*(s')) - Q^{\tilde{\pi}}(s', \tilde{\pi}(s'))]$$

$$\epsilon \leq 2\|Q^* - Q\|_{\infty} + \gamma \epsilon \sum_{s' \in S} p(s' | s, \tilde{\pi}(s))$$

i.e. $(1 - \gamma) \epsilon \leq 2\|Q^* - Q\|_{\infty}$

and

$$\boxed{\epsilon \leq \frac{2\|Q^* - Q\|_{\infty}}{1 - \gamma}}$$

