

Reinforcement Learning

Exercise 2 - Solution

Jonathan Schnitzler - st166934
Erick Villanueva Villaseñor - st190300
Eric Choquet - st160996

April 28, 2024

Proofs

a) Bellman optimality operator is a gamma-contraction We want to show

$$(\mathcal{T}v)(s) = \max_a \sum_{s',r} p(s', r|s, a)[r + \gamma v(s')] \quad (1)$$

fulfills the γ -contraction property, namely

$$\|\mathcal{T}v - \mathcal{T}w\|_\infty \leq \gamma \|v - w\|_\infty \quad (2)$$

Inspired by the lecture for the Bellman expectation backup operator, we will similarly use the definition of the infinity norm to show the contraction property

$$\|\mathcal{T}v - \mathcal{T}w\|_\infty = \left\| \max_a \sum_{s',r} p(s', r|s, a)[r + \gamma v(s')] - \max_a \sum_{s',r} p(s', r|s, a)[r + \gamma w(s')] \right\| \quad (3)$$

$$\leq \left\| \max_a \sum_{s',r} p(s', r|s, a)[r + \gamma v(s') - (r + \gamma w(s'))] \right\| \quad (4)$$

$$= \gamma \left\| \max_a \sum_{s',r} p(s', r|s, a)[v(s') - w(s')] \right\| \quad (5)$$

$$\leq \gamma \left\| \max_a \sum_{s',r} p(s', r|s, a) \|v(s') - w(s')\|_\infty \right\| \quad (6)$$

$$\leq \gamma \|v - w\|_\infty \quad (7)$$

b) Bounding general finite MDPs This is quite simple by imagining, a sequence of actions for which always the best reward r_{\max} or always the worst outcome, i.e. r_{\min} occurs. We can use the geometric sum formula for $\gamma < 1$

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] \quad (8)$$

$$= \mathbb{E}_\pi\left[\sum_{i=0}^{\infty} \gamma R_{t+i+1} | S_t = s\right] \quad (9)$$

$$\leq \sum_{i=0}^{\infty} \gamma r_{\max} \quad (10)$$

$$= r_{\max} \frac{1}{1 - \gamma} \quad (11)$$

which reversly holds for the minimum with a lower bound

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] \quad (12)$$

$$= \mathbb{E}_\pi\left[\sum_{i=0}^{\infty} \gamma R_{t+i+1} | S_t = s\right] \quad (13)$$

$$\geq \sum_{i=0}^{\infty} \gamma r_{\min} \quad (14)$$

$$= r_{\min} \frac{1}{1 - \gamma} \quad (15)$$

This yields

$$\frac{r_{\min}}{1 - \gamma} \leq v(s) \leq \frac{r_{\max}}{1 - \gamma} \quad (16)$$

From this we can follow from arbitrary $v(s)$ and $v(s')$ by assuming without loss of generality that $v(s) \geq v(s')$ (since the naming is arbitrary)

$$|v(s) - v(s')| = v(s) - v(s') \quad (17)$$

$$\leq \frac{r_{\max}}{1 - \gamma} - v(s') \quad (18)$$

$$\leq \frac{r_{\max}}{1 - \gamma} - \frac{r_{\min}}{1 - \gamma} \quad (19)$$

$$= \frac{r_{\max} - r_{\min}}{1 - \gamma} \quad (20)$$

which concludes the proof.

Value Iteration

a) Implementation of the value function

The value function is initialized with zero-values

$$V(s) = 0 \quad \forall s \in \mathcal{S} \quad (21)$$

and $\gamma = 0.8$, $\theta = 10^{-8}$.

0.498	0.832	1.311
0.536	0.977	2.295
0.306	0	5

Figure 1: Optimal value v_*

b) Optimal policy of value function