

# Reinforcement Learning

## Exercise 3 - Solution

Jonathan Schnitzler - st166934

Eric Choquet - st160996

April 29, 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  $\gamma$ -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

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

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

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

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

$$= \gamma \|v(s') - w(s')\|_\infty \left\| \max_a \sum_{s',r} p(s', r|s, a) \right\| \quad (7)$$

$$\leq \gamma \|v - w\|_\infty \quad (8)$$

**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  $\gamma < 1$

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] \quad (9)$$

$$= \mathbb{E}_\pi\left[\sum_{i=0}^{\infty} \gamma^i R_{t+i+1} | S_t = s\right] \quad (10)$$

$$\leq \sum_{i=0}^{\infty} \gamma^i r_{\max} \quad (11)$$

$$= r_{\max} \frac{1}{1 - \gamma} \quad (12)$$

which reversly holds for the minimum with a lower bound

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] \quad (13)$$

$$= \mathbb{E}_\pi\left[\sum_{i=0}^{\infty} \gamma^i R_{t+i+1} | S_t = s\right] \quad (14)$$

$$\geq \sum_{i=0}^{\infty} \gamma^i r_{\min} \quad (15)$$

$$= r_{\min} \frac{1}{1 - \gamma} \quad (16)$$

This yields

$$\frac{r_{\min}}{1 - \gamma} \leq v(s) \leq \frac{r_{\max}}{1 - \gamma} \quad (17)$$

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

$$|v(s) - v(s')| = v(s) - v(s') \quad (18)$$

$$\leq \frac{r_{\max}}{1 - \gamma} - v(s') \quad (19)$$

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

$$= \frac{r_{\max} - r_{\min}}{1 - \gamma} \quad (21)$$

which concludes the proof.

## Value Iteration

### a) Implementation of the value function

The value function is initialized with zero-values

$$V(s) = 0 \quad \forall_{s \in \mathcal{S}} \quad (22)$$

and  $\gamma = 0.8$ ,  $\theta = 10^{-8}$ . It converges in 43 Iterations

0.015	0.016	0.027	0.016
0.027	0.000	0.060	0.000
0.058	0.134	0.197	0.000
0.000	0.247	0.544	0.000

Table 1: Optimal value  $v_*$

$\downarrow$	$\uparrow$	$\rightarrow$	$\uparrow$
$\leftarrow$	H	$\leftarrow$	H
$\uparrow$	$\downarrow$	$\leftarrow$	H
H	$\rightarrow$	$\downarrow$	G

Table 2: Optimal policy  $\pi_*$

**b) Optimal policy of value function**