

Prof. Dr. Leif Döring  
Sara Klein, Benedikt Wille

Reinforcement Learning

## 9. Exercise Sheet

### 1. Sample based policy iteration without bounded rewards

Let the second moments of the rewards given a policy  $\pi \in \Pi_S$  exist, i.e.

$$\mathbb{E}_s^\pi [R_0^2] < \infty \quad \forall s \in \mathcal{S}.$$

Show that the Theorems 4.5.1 and 4.5.2 still apply to this policy, so that the one-step policy evaluation schemes from the lecture converge.

### 2. Convergence theorem 4.3.8 under weaker assumptions

Show that the statement of Theorem 4.3.8 also holds if  $\mathbb{E}[\varepsilon_i(n) | \mathcal{F}_n] \neq 0$  but instead satisfies

$$\sum_{n=1}^{\infty} \alpha_i(n) |\mathbb{E}[\varepsilon_i(n) | \mathcal{F}_n]| < \infty$$

almost surely for all coordinates  $i = 1, \dots, d$ . It is enough to prove an improved version of Lemma 4.4.4 where the condition  $\mathbb{E}[\varepsilon_n | \mathcal{F}_n] = 0$  is replaced with

$$\sum_{n=1}^{\infty} \alpha_n |\mathbb{E}[\varepsilon_n | \mathcal{F}_n]| < \infty. \tag{1}$$

Apply the Robbins-Siegmund theorem to  $W^2$  and use that  $W \leq 1 + W^2$ .

### 3. Programming task: One-step policy evaluation on grid world

We want to use the grid world example to illustrate how to perform policy evaluation:

- a) Implement the grid world example from the lecture notes with target in the lower right corner and trap diagonally above or modify the code from the lecture's webpage.
- b) Implement the Algorithms 17 and 18, the one-step policy evaluation schemes for  $V^\pi$  and  $Q^\pi$  respectively, for the grid world example.
- c) Think about what you intuitively think the best policy  $\pi^+$  and the worst policy  $\pi^-$  are for grid world and let additionally  $\pi$  be the policy that chooses the next action uniformly for all available options. Calculate  $n = 1000$  steps of each policy evaluation scheme for  $\pi^+$ ,  $\pi^-$ , and  $\pi$ .
- d) Compare Algorithm 17 to Algorithm 7, the iterative policy evaluation. Which algorithm do you think performs better? Can we always apply both algorithms?