

6. Exercise Sheet

1. Policy evaluation

Consider Algorithm 8 from the lecture. In Theorem 3.3.2 we proved convergence for this algorithm if $\gamma < 1$. Now assume $\gamma = 1$ and set $\Delta = 2\epsilon$ in the initialisation and choose termination condition $\Delta < \epsilon$. Give an example such that Algorithm 8 does not converge using $\gamma = 1$. Your are allowed to initialise the value function V arbitrarily.

2. Convergence of the in-place policy evaluation algorithm

Recall Algorithm 9 from the lecture. We aim to prove convergence of the algorithm (without termination) to V^π . Therefore label the state space \mathcal{S} by s_1, \dots, s_K and define

$$T_s^\pi V(s') = \begin{cases} T^\pi V(s) & : s = s' \\ V(s) & : s \neq s' \end{cases}$$

Define the composition $\bar{T}^\pi : U \rightarrow U$, $\bar{T}^\pi(v) := (T_{s_K}^\pi \circ \dots \circ T_{s_1}^\pi)(v)$ on the space of all functions $U = \{u : \mathcal{S} \rightarrow \mathbb{R}\}$ equipped with the supremums norm.

- Argue why \bar{T}^π is different from the Bellman operator T^π .
- Show that V^π is a fixpoint of the operator \bar{T}^π .
- Prove that \bar{T}^π is a contraction on $(U, \|\cdot\|_\infty)$.

3. Monte Carlo generalised ϵ -greedy policy iteration

In this exercise, we will use a Monte Carlo estimator to perform policy iteration. After Grid World has been a first example for a Markov decision process, we will now deal with more complex examples.

- Implement the ice vendor (Example 3.1.8) as a Markov decision process.
- Use the policy iteration from the last lecture to get the optimal decision rule for the iceman.
- Implement the Monte Carlo generalised ϵ -greedy policy iteration (algorithm 18) to get the optimal decision rule for the iceman. Additionally, use different ϵ parameters as well as a sequence $\epsilon_n \downarrow 0, n \rightarrow \infty$.