

Mathematical Foundations of Reinforcement Learning

 tuyaux.winak.be/index.php/Mathematical_Foundations_of_Reinforcement_Learning

Mathematical Foundations of Reinforcement Learning

Richting Informatica, Wiskunde

Jaar MWIS, MINE

Examenvragen

2022-2023 : Januari

Oefeningen

1. (5pt) Let $T=50$, and consider the following MDP (see image). Construct a stationary MDP with the same states and at most 2 actions such that all values are preserved.
2. (15pt)
 1. What is the period of the MDP in question 1?
 2. Give an example of a policy such that the resulting Markov chain has at least 2 different stationary distributions.
3. (10pt) Give an example of a stationary MDP such that the value iteration algorithm for discounted rewards converges in finitely many steps and one where it converges in infinitely many steps.
4. (10pt) Give an example of a stationary MDP such that from state i , no optimal unichain policy exists for the infinite-horizon expected average reward criterion.
5. (5 pt) For the Q-learning algorithm we defined learning rates $\gamma_0=0, \gamma_1, \gamma_2, \dots$ with $0 \leq \gamma_k < 1, \forall k \in \mathbb{N}$. Let $N_{i,a}$ be the map that maps n to the timestep at which the state-action pair (a,i) is visited for the n -th time. We required the following conditions: $\sum_{n=1}^{+\infty} \gamma_{N_{i,a}(n)} = +\infty$ and $\sum_{n=1}^{+\infty} \gamma_{N_{i,a}(n)}^2 < +\infty$. Give an example of learning rates that satisfy all the given conditions.

Theorie

1. (10pt) Let M be a stationary communicating MDP. Prove that there is an optimal unichain policy from i for the infinite-horizon expected limit-average reward criterion.

2. (10pt)

1. Consider the first condition of question 5 of the exercises. Argue that this implies that every state-action pair is seen infinitely often with probability 1.
2. Give an example of a communicating MDP, initial values $Q(0)$ such that even under the assumptions, some state-action pair may be visited finitely often with probability 0.

3. (15pt) We will now focus on the two mappings $U: \mathbb{R}^N \rightarrow \mathbb{R}^N$ and $L_\pi: \mathbb{R}^N \rightarrow \mathbb{R}^N$ defined for any randomized decision rule π as follows.

$\{Ux\}_i = \max_{a \in A(i)} \{r_i(a) + \alpha \sum_{j \in S} p_{ij}(a)x_j\}$ and $\{L_\pi x\}_i = r_i(\pi) + \alpha \sum_{j \in S} p_{ij}(\pi)x_j$. And

$\{L_\pi x\}_i = r_i(\pi) + \alpha \sum_{j \in S} p_{ij}(\pi)x_j$. Prove the following result. (You must recall the definition of monotone contraction mapping. Also, you may use results proved during the course as auxiliary lemmas. "The mappings U and L_π for any randomized decision rule π are monotone contraction mappings with contraction factor α . Moreover, $v_\alpha(\pi^\infty)$ is the unique fixed point of $L_\pi x = x$."

Moreover, $v_\alpha(\pi^\infty)$ is the unique fixed point of $L_\pi x = x$."

4. (10pt) For the policy iteration algorithm for infinite-horizon discounted rewards, we defined the action set $A(i, f)$ as follows:

$A(i, f) := \{a \in A(i) \mid r_i(a) + \alpha \sum_{j \in S} p_{ij}(a)v_j(f) > v_i(f)\}$

Take $i \in S$ and $f \in C(D)$. Prove the following statements:

1. If $A(i, f) = \emptyset$, for every $i \in S$, then f is an α -discounted optimal policy.
2. If $A(i, f) \neq \emptyset$, for some $i \in S$, then $v_\alpha(g) > v_\alpha(f)$ for any $g \in C(D)$ with $g(i) \in A(i, f)$ when $g(i) \neq f(i)$.