

Abstract

abstract in German in at least 500 words is a requirement!!!

optionally also in English

write this last with introduction

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 2 | Background on Formal Languages | 2 |
| 3 | Background on Reinforcement Learning | 3 |
| 4 | Reward Machines | 7 |
| 5 | Conclusion & Outlook | 8 |
| | References | 9 |
| A | Appendix | 10 |

1 Introduction

write this last with abstract

2 Background on Formal Languages

write after two other chapters (15-20 pages)

preliminaries to cover:

- regular languages (LTL, regular expressions)
- (deterministic) finite state automata (DFA)
- mealy/moore machines
- büchli automata

3 Background on Reinforcement Learning

(10-15 pages)

cover:

- MDP
- policy
- q function
- Bellman equations
- q learning

The goal of reinforcement learning is to learn a task by maximizing the total rewards along sequences of decisions called episodes. The difficulty is that decisions can have delayed consequences.

too terse, explain more naturally

We have an agent that acts in an environment. How the environment behaves is unknown to the agent. Furthermore for each action the agent takes it receives a reward from the environment, but how these rewards are determined is also not visible to the agent. So the environment is a complete black box to the agent.

Consider the following example as a demonstration of the reinforcement learning framework. In a grid-based office environment a robot is tasked with bringing coffee from the kitchen to a particular office, see figure 1 for the layout.

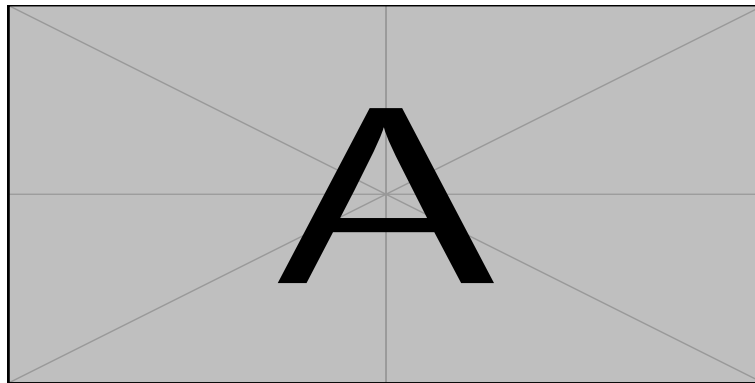


Figure 1: office world...

The environment starts in state S_0 and for each state $S_t \in \mathcal{S}$ at time $t \in \mathbb{N}$ the agent has to decide on an action $A_t \in \mathcal{A}$ like going up, down, left, right or use the coffee machine. After executing the action the environment changes to a new state S_{t+1} and rewards the robot with a reward $R_{t+1} \in \mathcal{R} \subseteq \mathbb{R}$. This feedback loop is shown as a diagram in figure 2.

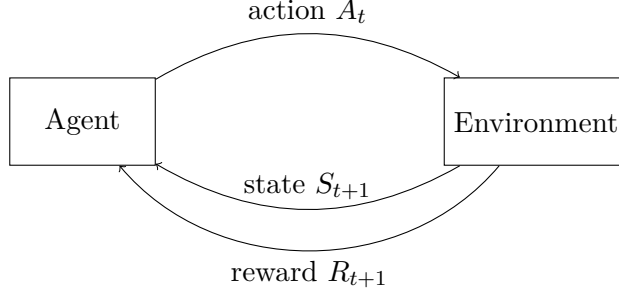


Figure 2: Diagram of the reinforcement learning feedback loop.

From the point of view of the agent the environment and the rewards are random processes. Here for example $(S_t)_{t \in \mathbb{N}} \subseteq \mathcal{S}$ describes the evolution of the states and $(R_t)_{t \in \mathbb{N}} \subseteq \mathcal{R}$ the rewards, which are defined on some probability space with joint measure \mathbb{P} . A key assumption is that these processes do not remember any history other than the state that they are currently in, i.e. $\mathbb{P}(S_{t+1} = s, R_{t+1} = r \mid S_0, A_0, S_1, R_1, A_1, \dots, S_t, R_t, A_t) = \mathbb{P}(S_{t+1} = s, R_{t+1} = r \mid S_t, A_t)$. This means that we assume that the environment has the Markov property and although in real life scenarios this might not be quite true it is still a good approximation and the theory for non-Markovian processes also uses the ideas we will develop in this chapter.

Following [1] we will also assume finite spaces for the sake of clarity but the theory can also be developed using probability densities on infinite spaces. Usually we do not care about the details of these processes and underlying probability spaces, but we work with some derived quantities instead, namely the state-transition probability distribution $p(s' | s, a) = \sum_{r \in \mathcal{R}} \mathbb{P}(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a)$ and the expected reward $r(s, a) = \mathbb{E}(R_{t+1} \mid S_t = s, A_t = a) = \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} \mathbb{P}(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a)$. This motivates the defining model of RL, the Markov decision process:

why markov property? is there some more motivation other than it works?

Definition 3.1. A *Markov Decision Process (MDP)* is given by a tuple $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$, where \mathcal{S} is a finite set of *states*, \mathcal{A} is a finite set of *actions*, $p : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the *state-transition probability distribution*, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the *reward function* and $\gamma \in [0, 1]$ is a *discount factor*. Note that here implicitly \mathcal{R} is the image of r and thus a finite subset of \mathbb{R} .

This MDP models the diagram in figure 2 from above as follows: From a given state $s \in \mathcal{S}$ the agent chooses an action $a \in \mathcal{A}$ and the environment changes to state $s' \in \mathcal{S}$ with probability $p(s' | s, a)$ and gives reward $r(s, a)$. The discount factor will become relevant in a moment but in essence a lower discount factor would motivate the agent to take actions based on the reward sooner rather than later as with a higher discount factor.

In the office world example ... TODO: what is the mdp in the example

TODO: motivate a policy

Definition 3.2. A policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is a probability distribution that prescribes an action $a \in \mathcal{A}$ to be taken for a given state $s \in \mathcal{S}$ with probability $\pi(a|s)$.

For our office world example a policy could be to go left when in room 1 and go right when in room 2 or stand still when not in either room. Of course this would not be a good policy, but what constitutes a "good" policy will be defined with the state value function and the action value function.

Definition 3.3. The state value function, often just value function, $v_\pi : \mathcal{S} \rightarrow \mathbb{R}$ under a policy π is defined as the expected discounted sum of future rewards starting from a state $s \in \mathcal{S}$ and taking actions according to π :

$$v_\pi(s) = r_\pi(s) + \gamma \sum_{s' \in \mathcal{S}} \left(p_\pi(s'|s) r_\pi(s') + \gamma \sum_{s'' \in \mathcal{S}} \left(p_\pi(s''|s') r_\pi(s'') + \dots \right) \right)$$

where $p_\pi(s'|s) = \sum_{a \in \mathcal{A}} \pi(a|s) p(s'|s, a)$ and $r_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) r(s, a)$. Since we only consider finite spaces in this setting we get the following in vector form:

$$\begin{aligned} V_\pi &= R_\pi + \gamma \sum_{s' \in \mathcal{S}} \left((P_\pi)_{s,s'} (R_\pi)_s + \gamma \sum_{s'' \in \mathcal{S}} \left((P_\pi)_{s',s''} (R_\pi)_{s''} + \dots \right) \right) \\ &= R_\pi + \gamma \sum_{s' \in \mathcal{S}} (P_\pi)_{s,s'} (R_\pi)_{s'} + \gamma^2 \sum_{s' \in \mathcal{S}} (P_\pi^2)_{s,s'} (R_\pi)_{s'} + \dots \\ &= R_\pi + \gamma P_\pi V_\pi \end{aligned}$$

where $P_\pi \in [0, 1]^{|\mathcal{S}| \times |\mathcal{S}|}$ is the transition matrix, $R_\pi \in \mathbb{R}^{|\mathcal{S}|}$ the vector of all rewards and $V_\pi \in \mathbb{R}^{|\mathcal{S}|}$ the state-values vector.

Definition 3.4. The action value function or Q-function $q_\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ under a policy π is the expected discounted sum of taking action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$ and then following policy π :

$$q_\pi(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_\pi(s')$$

Now that we have defined the value of a policy for each state we can define when a policy is considered "better" than another.

Definition 3.5. We say $\pi \geq \pi'$ if $v_\pi(s) \geq v_{\pi'}(s)$ for all states $s \in \mathcal{S}$, which naturally defines a partial order on the space of policies.

Of course this leads to the natural question if there is a "best" policy, to which the answer is yes.

Proposition 3.6. For a finite MDP there exists an optimal policy.

Proof. TODO!

□

write about the dynamical programming approach: https://www.wikiwand.com/en/Dynamic_programming)

4 Reward Machines

(15-20 pages)

papers:

- original reward machine paper from 2018 [2]
- connection to LTL from 2019 [3]
- newer reward machine paper from 2022 [4]

include updated RL framework graph to better differentiate that reward machines are not in the environment black box

include same task from previous chapter and how a reward machine for that might look like as a graph

central theorems:

- construction/correspondence to LTL and such
- convergence proof

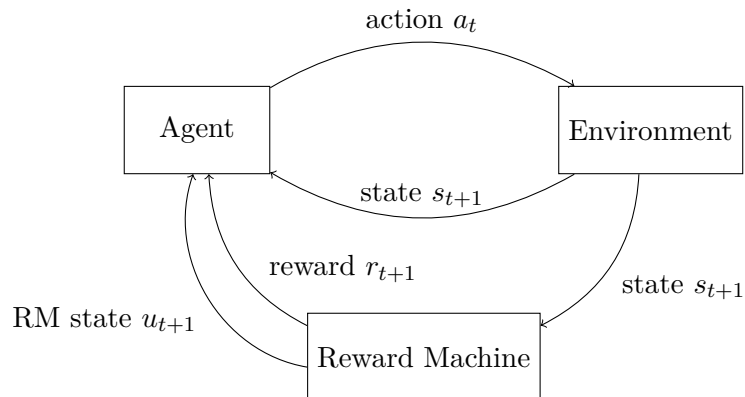


Figure 3: Diagram of the reinforcement learning feedback loop with a reward machine. The agent now not only gets an environment state but also a reward machine state.

The mathematical model for this is ...

Definition 4.1. A *Markov Decision Process with Reward Machine (MDPRM)* is ...

5 Conclusion & Outlook

write this after the main chapters

References

- [1] R.S. Sutton and A.G. Barto. *Reinforcement Learning, second edition: An Introduction*. Adaptive Computation and Machine Learning series. MIT Press, 2018. ISBN: 9780262039246. URL: <https://books.google.at/books?id=5s-MEAAAQBAJ>.
- [2] Rodrigo Toro Icarte et al. “Using Reward Machines for High-Level Task Specification and Decomposition in Reinforcement Learning”. In: *Proceedings of the 35th International Conference on Machine Learning (ICML)*. 2018, pp. 2112–2121. URL: <http://proceedings.mlr.press/v80/icarte18a.html>.
- [3] Alberto Camacho et al. “LTL and Beyond: Formal Languages for Reward Function Specification in Reinforcement Learning”. In: *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI)*. 2019, pp. 6065–6073. URL: <https://www.ijcai.org/proceedings/2019/0840.pdf>.
- [4] Rodrigo Toro Icarte et al. “Reward Machines: Exploiting Reward Function Structure in Reinforcement Learning”. In: *Journal of Artificial Intelligence Research (JAIR)* 73 (2022), pp. 173–208. URL: <https://doi.org/10.1613/jair.1.12440>.

A Appendix

experiments and code go here convergence graphs of qrm/crm and standard