

Reinforcement learning

Jiří Kléma

Department of Computer Science,
Czech Technical University in Prague

Lecture based on **AIMA book** and its accompanying slides



<https://cw.fel.cvut.cz/wiki/courses/b4m36smu>

Reinforcement learning (RL) in a nutshell

Imagine playing a new game whose rules you don't know;
after a hundred or so moves, your opponent announces, "You lose".

Russell and Norvig
Introduction to Artificial Intelligence

Agenda

- A formal definition of RL
 - the distinction from MDP known from ZUI course,
 - a definition in terms of our SMU general framework,
 - (necessary) simplifications from the most general form.
- Passive learning
 - model-free and model-based approaches,
 - direct utility estimation, adaptive dynamic programming (ADP), temporal difference (TD) learning,
 - illustrative, but not really useful in practice.
- Active learning
 - exploration-exploitation dilemma,
 - active ADP agent, Q-learner based on TD learning.
- Generalizations and applications
 - utility function approximations in large state spaces.

Passive and active learning

■ Passive learning

- the agent employs a fixed policy π ,
- it learns how good the policy is by interaction with the environment, e.g., it learns $U^\pi(s)$,
- analogous to policy evaluation in policy iteration, (however, the big difference is that the model of the environment is unknown now).

■ Active learning

- the agent searches for an optimal (or at least good) policy,
- it **explores** (many) different actions in (many) different states,
- analogous to learning and solving the underlying MDP.

Model-based and model-free learning

- Model-based learning

- learn the MDP model (\mathcal{S} and R), or an approximation of it,
- use the model to guess the state utility and find optimal policy,
- more difficult to tailor to the environment
(as straightforward application can be slow).

- Model-free learning

- derive the optimal policy without explicitly learning the environment model,
- typically stem from the estimation of Q , the utility of state-action pairs,
- easier to apply independently of the environment,
may fail in complex worlds.

