

Symbolic Model-Based Reinforcement Learning

Simone Manti, AU ID 734894, smanti@mpe.au.dk
Final project for Neural Networks and Deep Learning 2024

October 12, 2024

1 Introduction

2 Model-Based Reinforcement Learning

Reinforcement Learning (RL), contrarily to standard Supervised Learning, requires a continue interaction with data to update the policy. Usually, model-free RL is employed to optimize the policy, where the agent can only use the data sampled interacting with the real environment. In particular, RL algorithms with high sample complexity are non-manageable in a real-world scenario due to its computational cost. An alternative is presented by the so-called Model-Based Reinforcement Learning (MBRL) [2], where apart from optimizing policy the environment model must be learned. In real-world scenario, having an environment model drastically reduces the complexity of generating new training samples for RL algorithms. In the following, we briefly overview the main components of MBRL. A more extensive treatment can be found in [2, 3, 4].

MBRL restates the problem as a Markov Decision Process (MDP) [1]. A MDP is a tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, p(s_0), \gamma\}$, where \mathcal{S} and \mathcal{A} are, respectively, the state and the action space, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow p(\mathcal{S})$ is the *transition function*, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the *reward function*, $p(s_0)$ is the initial state distribution and $\gamma \in [0, 1]$ is a discount parameter. A solution to a finite-time MDP consists of a policy $\pi^* : \mathcal{S} \rightarrow \mathcal{A}$ which maximizes the expected sum of rewards, that is

$$\pi^* = \arg \max_{\pi} Q^{\pi}(s, a) = \arg \max_{\pi} \mathbb{E}_{\pi, \mathcal{T}} \left[\sum_{k=0}^K \gamma^k r_{t+k} | s_t = s, a_t = a \right], \quad (1)$$

where $s_t \in \mathcal{S}, a_t \in \mathcal{A}, r_t := \mathcal{R}(s_t, a_t, s_{t+1})$.

MBRL approaches repeat the following two steps iteratively:

1. collect (or enlarge the current data with) a dataset $\mathcal{D} := \{(s_i, a_i, s_{i+1}, r_i, d_i)\}_{i=1}^N$, where d_i is the termination value indicating the end of the episode, and learn an approximate model f^* of the environment's dynamics fitting \mathcal{D} , i.e.

$$f^* = \arg \max_{f \in \mathcal{F}} \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}} \mathcal{L}(s_{t+1}, f(s_t, a_t)), \quad (2)$$

where \mathcal{F} is a general family of functions (for instance Neural Networks) and \mathcal{L} is a loss function. Note that with the same data reward and termination functions can be learned reward as well;

2. simulate the transition with the learned model f^* and optimize the policy maximizing the expected sum of rewards like in (1).

3 Symbolic Regression

SR consists of finding a mathematical expressions that best fits a given dataset without assuming a particular structure of the expression. In our case, SR is handled with Genetic Programming (GP) [5, 6]. GP is an evolutionary strategy that explores a space of candidate models/functions by iteratively updating an initial population through genetic operations. A candidate model is also called *individual* and can be easily represented by a tree, in which each node is either a *primitive* (i.e. an operator) or a *terminal* (i.e. a variable or constant).

Initialization The population is typically initialized randomly. Some well-known techniques are the *full* and the *grow* method (see Section 2.1 in [6]), and a combination of those gives rise to *ramped half and half*. Each of them assumes that the starting population contains individuals that do not exceed a fixed maximum depth.

Selection Usually genetic operations are performed on individuals that are selected based on fitness, i.e. that are “probably good”. One way to do so is *tournament selection*. This consists of performing a tournament for each k (tournament size) individuals, competing in terms of their fitnesses. The higher the score in these tournaments, the higher the probability for these individuals to be chosen. In our method, $k = 2$ and we assign a non-zero probability to pick the loser instead of the winner in each tournament. This is calibrated in the hyperparameter selection phase.

Crossover and Mutation Genetic operations consists of *crossover* and *mutation*. The most common crossover used is the *subtree crossover* (or *one-point crossover*): given two parent trees, a random subtree and a crossover point are chosen for both; then, two new trees are generated, named *offsprings*, replacing the subtree rooted in the first parent with the one in the second parent and vice-versa. Instead, the most common mutation creates an offspring for a given parent tree by replacing a (randomly chosen) subtree with a newly generated subtree.

4 Experiments and implementations details

Add intro + libraries used

4.1 Simple 1-dimensional MDP

The following simple 1-dimensional MDP, with episode length 10, is considered

$$s_{t+1} := s_t + a_t, \quad r_t := \cos(2\pi s_{t+1}) \exp(|s_{t+1}|/3). \quad (3)$$

In this case, the goal is learning both the dynamics and the reward function. [Continue from here](#)

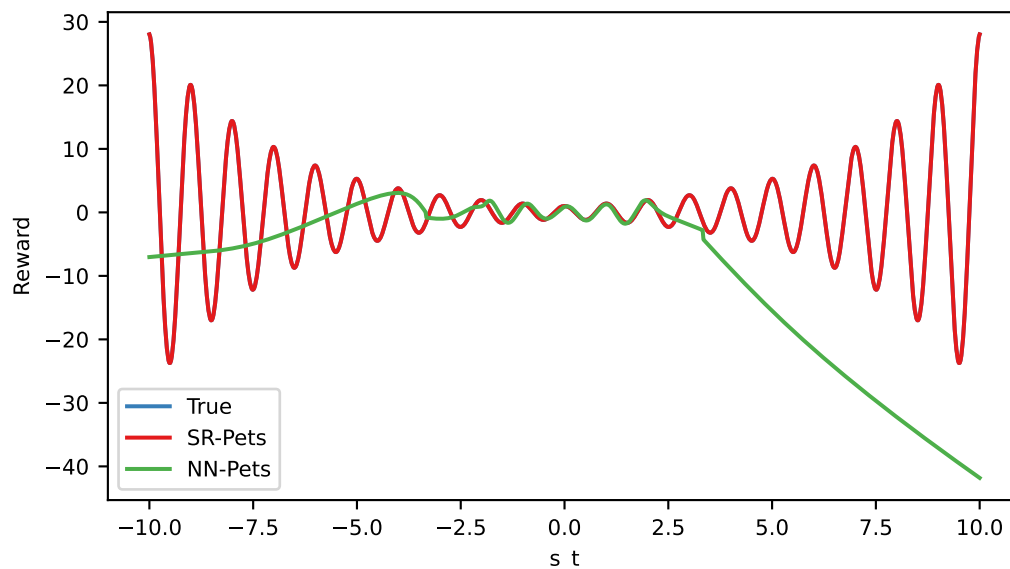


Figure 1: ??

4.2 Cartpole

5 Conclusion and discussion

References

- [1] M. L. Puterman, “Markov Decision Processes: Discrete Stochastic Dynamic Programming.” *John Wiley and sons*, 2014.
- [2] T.M. Moerland, J. Broekens, A. Plaat, C.M. Jonker, “Model-based Reinforcement Learning: A Survey”. *Foundations and Trends in Machine Learning*, 2023.
- [3] L. Pineda, B. Amos, A. Zhang, N.O Lambert, R. Calandra, “Mbrl-lib: A Modular Library for Model-Based Reinforcement Learning”. *arXiv preprint arXiv:2104.10159*, 2021.
- [4] F. Luo, T. Xu, H. Lai, X. Chen, W. Zhang, Y. Yu, “A Survey on Model-Based Reinforcement Learning”. *Science China Information Sciences*, 2024.
- [5] J.R. Koza, “Genetic programming as a means for programming computers by natural selection”, *Statistics and Computing*, vol. 4, 1994.
- [6] R. Poli et al., “A Field Guide to Genetic Programming”. Published via <http://lulu.com> and freely available at <http://www.gp-field-guide.org.uk>, 2008.