

When to use parametric models in reinforcement learning?

Growing advances in model-based reinforcement learning have yielded data-efficient methods. These methods learn a model of the dynamics of the world and align behaviors of the agent with the model's beliefs. However, it is often unclear as to when to use a model for acting and planning since models may inherently be imperfect in nature and can steer the agent towards a sub-optimal policy in the case of inaccurate beliefs. The work investigates the usage of parameterized models and their relationship to replay memory in various reinforcement learning settings. Replay memory is thought of as analogous to a model being used for behavioral updates. Such behavioral models are used for training agents in a sample-efficient manner in comparison to pure model-based approaches.

Parameterized models serve as an internal component for the agent in order to plan for future steps. However, planning is hindered by inaccurate beliefs in the long-horizon and computational expense since the model is updated to learn complete dynamics of the environment. Pure model-based approaches collect data from the environment and train the agent based on the model being updated by the agent's policy. Another method for using models is by executing inverse updates which predicts backwards. Learning an inverse model is beneficial since it only provides fictional states and does not harm agent's behaviors. On the other hand, a replay-based approach collects data from the environment and simply updates the agent's policy. A stark contrast in learning with these methods can be observed on the basis of (1) scalability and (2) performance. Firstly, the conventional planning approach is scalable in the number of planning steps and reduces data dependency as the planning horizon is improved. Secondly, inverse model learning is comparable replay in terms of performance under high stochasticity. This indicates that a pure model-based approach is not robust to changes in environment dynamics. The claim is further strengthened by carrying out a large-scale study of the novel data-efficient Rainbow which uses fewer samples to demonstrate better performance in comparison to state-of-the-art Simulated Policy Learning (SimPLE) method.