

Model Predictive Control and Reinforcement Learning – Planning and Learning –

Joschka Boedecker and Moritz Diehl

University Freiburg

July 30, 2021





1 Model Learning

2 Dyna

3 Simulation-based Search

Acknowledgement



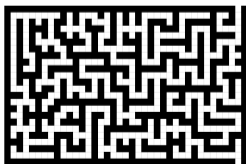
Slide contents are partially based on *Reinforcement Learning: An Introduction* by Sutton and Barto and the Reinforcement Learning lecture by David Silver.



- ▶ Policy: defines the behaviour of the agent
 - ▶ is a mapping from a state to an action
 - ▶ can be stochastic: $\pi(a|s) = \mathbb{P}[A_t = a|S_t = s]$
 - ▶ or deterministic: $\pi(s) = a$
- ▶ Value-function: defines the expected value of a state or an action
 - ▶ $v_\pi(s) = \mathbb{E}[G_t|S_t = s]$ and $q_\pi(s, a) = \mathbb{E}[G_t|S_t = s, A_t = a]$
 - ▶ can be used to evaluate states or to extract a good policy
- ▶ Model: defines the transitions between states in an environment
 - ▶ p yields the next state and reward
 - ▶ $p(s', r|s, a) = \Pr\{S_{t+1} = s', R_{t+1} = r|S_t = s, A_t = a\}$

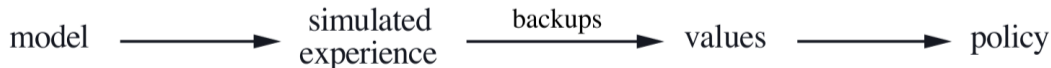
Learning Models

- ▶ Depending on the task, the dynamics model can be much easier than the value-function or the policy
- ▶ We can estimate it via supervised learning methods
- ▶ However, the model can also be more complex than policy and value-function
- ▶ In practice, modelling state-changes can even be easier than the global state
- ▶ In a nutshell:
 - ▶ Learning a model: data-efficient, hard to extract an optimal policy
 - ▶ Learning a value function: less data-efficient, easier to extract an optimal policy
 - ▶ Learning a policy: data-inefficient, directly estimate an optimal policy





- ▶ Given a state and an action, a model generates the next state and the corresponding reward (can also be used to generate sequences of states and rewards)
- ▶ It can either give the probabilities of all possible next states and rewards (*distribution model*), or only one (*sample model*)
- ▶ Which one was used in Dynamic Programming?
- ▶ Extracting a policy from a model is called *planning*





- ▶ Planning: Uses simulated experience generated by a model
- ▶ Learning: Uses real experiences from the environment
- ▶ But we can also apply learning methods to simulated experience

Random-sample one-step tabular Q-planning

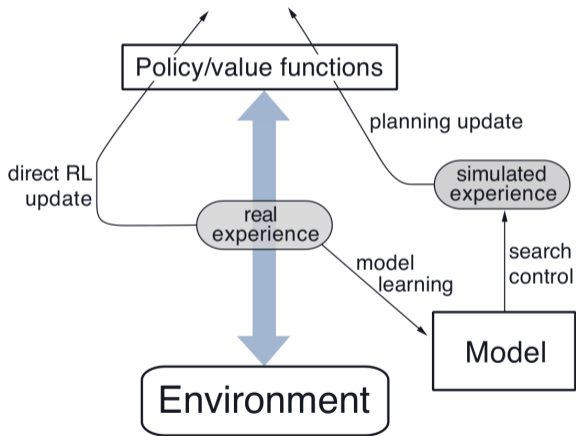
Loop forever:

1. Select a state, $S \in \mathcal{S}$, and an action, $A \in \mathcal{A}(S)$, at random
2. Send S, A to a sample model, and obtain
a sample next reward, R , and a sample next state, S'
3. Apply one-step tabular Q-learning to S, A, R, S' :
$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

- ▶ Converges to the optimal policy *for the model*



- ▶ Real experience can be used to optimize the value function (or the policy)
 - ▶ *directly* (model-free RL) or
 - ▶ *indirectly* (model-based RL) via a model
- ▶ *Indirect methods* are often more data-efficient
- ▶ But they introduce additional bias through the model
- ▶ Idea of Dyna: try to combine the best of both worlds



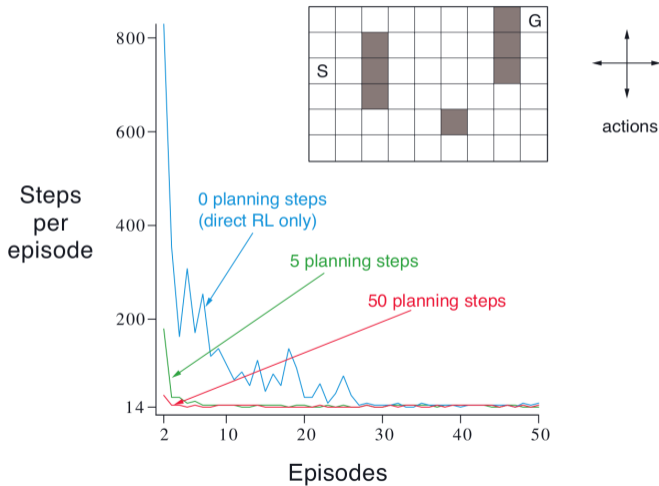


Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \varepsilon$ -greedy(S, Q)
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- (f) Loop repeat n times:
 - $S \leftarrow$ random previously observed state
 - $A \leftarrow$ random action previously taken in S
 - $R, S' \leftarrow Model(S, A)$
 - $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$





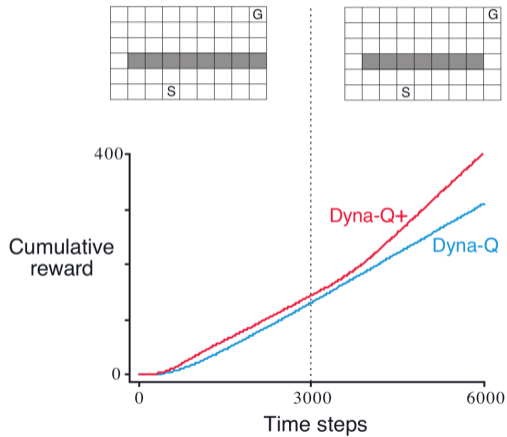
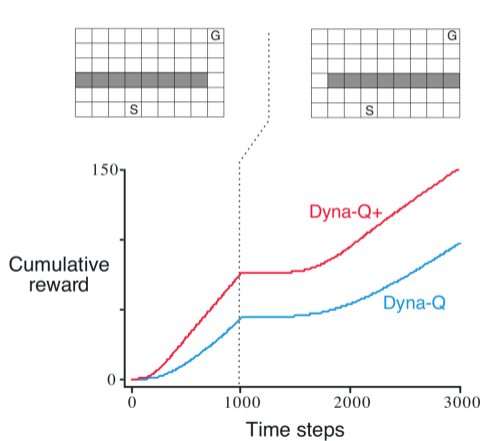
- ▶ Models can be incorrect (limited number of samples, environment has changed, function approximation)
- ▶ Especially in areas where the agent has not explored
- ▶ There can be a *Distribution Mismatch* when the agent enters new areas of the state-action space
- ▶ When the model is incorrect, the planning process is likely to find a suboptimal policy



- ▶ Dyna-Q+: Add an exploration bonus for transitions that have not been visited recently
- ▶ Let r be the reward, κ the weight of the exploration bonus and τ the number of time steps in which a certain transition has not been visited
- ▶ Then Dyna-Q+ modifies the internal reward function to:

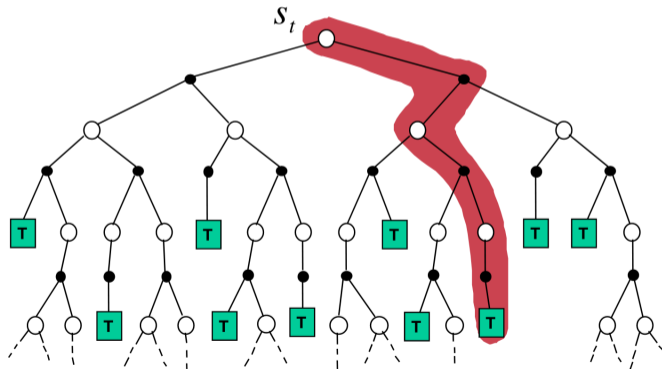
$$r + \kappa\sqrt{\tau}$$

When the model is wrong



Simulation-based Search

- ▶ Forward search paradigm using sample-based planning
- ▶ Simulate episodes of experience from now with the model
- ▶ Apply model-free RL to simulated episodes





- ▶ Build a search tree containing visited states and actions using the model (simulate episodes from current state)
- ▶ Two policies: tree policy (improving, e.g. ϵ -greedy) and out-of-tree rollout policy (random)
- ▶ Monte-Carlo control applied to simulated experience
- ▶ One of the key ingredients of AlphaGo (2016)



1. Selection: starting at the root, traverse to a leaf node following the *tree policy*
2. Expansion: expand the tree by one or multiple child nodes reached from the selected leaf in some iterations
3. Simulation: simulate an episode following the *rollout policy*
4. Backup: update the action-values for all nodes visited **in the tree**

Monte Carlo Tree Search

