Reinforcement Learning

Lecture 9: Model-based methods

Chris G. Willcocks Durham University



Lecture covers chapter 8 in Sutton & Barto [1] and examples from David Silver [2]

1 Model-based reinforcement learning

- taxonomy
- overview
- the simulation cycle
- characteristics

2 Integrated learning and planning

- Dyna-Q
- characteristics
- Monte Carlo tree search
- simulated policy learning

Model-based RL taxonomy



This figure does not capture overlap, for example between policy optimsiation and Q-learning algorithms

Model-based RL overview





In **model-free** RL:

- No model
- **Learn** the value function q(s, a) and/or the policy $\pi(a|s)$ from experience

In model-based RL:

- Learn the model from experience
- **Plan** the value function and/or the policy from the model





Model-based RL cycle:

- The agent experiences the real environment
- We learn a model to predict what the real environment does (when you take an action)
- We then use this simulated model to plan
- This allows us to estimate the value function and/or policy without directly interacting with the real environment
- But we use this policy to take real actions again



Model-based RL advantages:

- The model can sometimes be a simpler and more useful representation of the environment than you can otherwise access by experience
- Can be learnt by supervised learning
- Can reason about model uncertainty

Model-based RL disadvantages:

- This is another component which introduces some approximation error
 - Value function and/or policy approximation and now model approximation
- We can only be as good as our model





Definition: model

A model $\mathcal{M} = \langle \mathcal{P}_{\eta}, \mathcal{R}_{\eta} \rangle$ is a parameterised η representation of an MDP: $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R} \rangle$. It approximates state transitions $\mathcal{P}_{\eta} \approx \mathcal{P}$ and rewards $\mathcal{R}_{\eta} \approx \mathcal{R}$, learning a distribution over the next states and rewards:

 $S_{t+1} \sim \mathcal{P}(S_{t+1}|S_t, A_t)$ $R_{t+1} = \mathcal{R}(R_{t+1}|S_t, A_t),$

which typically are conditionally independent of each other:

$$P(S_{t+1}, R_{t+1}|S_t, A_t) = P(S_{t+1}|S_t, A_t)P(R_{t+1}|S_t, A_t)$$

Example: environment model



Learning the model

We learn the model M_{η} from experience $\{S_1, A_1, R_2, ..., S_T\}$ using **supervised learning**.

- We receive a stream of actual experiences
- This gives us a dataset:

 $S_1, A_1 \to R_2, S_2$ $S_2, A_2 \to R_3, S_3$

• $s, a \rightarrow r$ is a regression problem

• $s, a \rightarrow s'$ is a density estimation problem

Experience can be simulated and real

Simulated experience sampled from \mathcal{M}_{η}

 $S' \sim \mathcal{P}_{\eta}(S'|S, A)$ $R = \mathcal{R}_{\eta}(R|S, A)$

Real experience sampled from the true MDP

$$S' \sim \mathcal{P}^a_{s,s'}$$
$$R = \mathcal{R}^a_s$$





Algorithm: Dyna-Q [3, 4]

```
 \begin{array}{ll} \mbox{initialise } Q(s,a) \mbox{ and model } \mathcal{M}(s,a) \mbox{ for all } s \in \mathcal{S} \mbox{ and } a \in \mathcal{A}(s) \\ \mbox{while True:} \\ s \leftarrow \mbox{current (nonterminal) state} \\ a \leftarrow \epsilon - \mbox{greedy}(s,Q) \\ r,s' \leftarrow \mbox{env.step}(s,a) \\ Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \max_{\hat{a}} Q(s',\hat{a}) - Q(s,a)) \\ \mathcal{M}(s,a) \leftarrow r,s' \mbox{ (assuming deterministic environment)} \\ \mbox{for i in range(n):} \\ s \leftarrow \mbox{ random previously observed state} \\ a \leftarrow \mbox{ random action previously taken in } s \\ r,s' \leftarrow \mathcal{M}(s,a) \\ Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \max_{\hat{a}} Q(s',\hat{a}) - Q(s,a)) \\ \end{array}
```

Model-based RL Dyna-Q characteristics





"Model-Based Reinforcement Learning for Atari" [5]





Summary

In summary, model-based methods:

- are easy to train with supervised learning
- allow for planning ahead
- can be very data efficient
- can be used to imagine situations without experiencing them
- but the value and policy learnt can only be as good as the model
- they can be combined with model-free methods

References I

- [1] Richard S Sutton and Andrew G Barto.
 <u>Reinforcement learning: An introduction (second edition)</u>. <u>Available online</u> . MIT press, 2018.
- [2] David Silver. Reinforcement Learning lectures. https://www.davidsilver.uk/teaching/. 2015.
- [3] Richard S Sutton. "Integrated architectures for learning, planning, and reacting based on approximating dynamic programming". In: Machine learning proceedings 1990. Elsevier, 1990, pp. 216–224.
- [4] Baolin Peng et al. "Deep Dyna-Q: Integrating planning for task-completion dialogue policy learning". In: <u>arXiv preprint arXiv:1801.06176</u> (2018).
- [5] Lukasz Kaiser et al. "Model-based reinforcement learning for atari". In: arXiv preprint arXiv:1903.00374 (2019).