

# Reinforcement Learning Tutorial 7, Week 8

## Reward expectation, Prioritisation, and Hyperparameters\*

Pavlos Andreadis, Michael Herrmann

March 2025

### Problem 1 – $R$ -learning

$R$ -learning<sup>1</sup> is similar to  $Q$ -learning, in particular for non-discounted, non-episodic problems. It is based on the average reward  $\rho = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n E[r_t]$ , and considers the current rewards in comparison to this accumulating reward average towards the value:

$$V(s_t) = \sum_{k=1}^{\infty} E[r_{t+k} - \rho | s_t = s]$$

$$Q(s_t, a_t) = \sum_{k=1}^{\infty} E[r_{t+k} - \rho | s_t = s, a_t = a]$$

In this relative value function (relative to the average),  $\rho$  is slowly adapted as a measure of success. In this way a different concept of optimality is implied in particular for non-episodic tasks. As an algorithm,  $R$ -learning works as follows

1. Initialise  $\rho$  and  $Q(s, a)$
2. Observe  $s_t$  and choose  $a_t$  (e.g.  $\epsilon$ -greedy), execute  $a_t$
3. Observe  $r_{t+1}$  and  $s_{t+1}$
4. Update

$$Q_{t+1}(s_t, a_t) = (1 - \eta) Q_t(s_t, a_t) + \eta (r_{t+1} - \rho_t + \max_a Q_t(s_{t+1}, a))$$

5. If  $Q(s_t, a_t) = \max_a Q(s_t, a)$  then

$$\rho_{t+1} = (1 - \alpha) \rho_t + \alpha (r_{t+1} + \max_a Q_t(s_{t+1}, a) - \max_a Q_{t+1}(s_t, a))$$

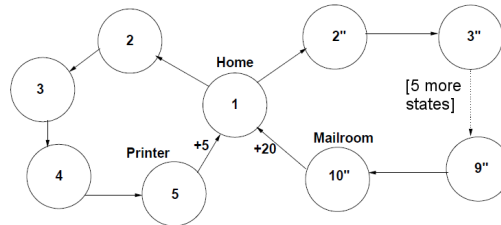
---

\*with special thanks to Adam Jelley

<sup>1</sup>A. Schwartz (1993) A reinforcement learning method for maximizing undiscounted rewards. 10th ICML. (You don't need to know  $R$ -learning for the exam.)

[We would choose  $\eta \gg \alpha$ , because, otherwise, for  $r = 0$ ,  $Q$ -value may cease to change and the agent may get trapped in a suboptimal limit cycle.]

Compare  $R$ -learning and  $Q$ -learning in the following simple example, where only one decision needs to be taken: The agent moves either to nearby printer (“o.k.”) or to distant mail room (“good”).



Is it possible that  $R$  learning finds the optimal solution quicker than  $Q$ -learning? How does the result for  $Q$ -learning depend on the discount factor  $\gamma$ ?

## Problem 2 – Prioritised sweeping<sup>2</sup>

While Dyna<sup>3</sup> agents select state-action pairs uniformly at random from the previously experienced pairs, it might be more efficient to use a non-uniform probability distribution. Why? Which state-action pairs should be preferred? Discuss the role of goal states in this context.

## Problem 3 - Discussion

### Part a

Considering a Reinforcement Learning algorithm in general, what is the overall effect of increasing the learning rate? What happens when you set it too high? What happens when you set it too low?

### Part b

Is the discount factor  $\gamma$ :

1. Part of the definition of a Markov Decision Process? That is, a part of the definition of the problem to be solved; or
2. Is it external to the problem? That is, a hyperparameter for training the model.

When a discount factor is close to 1, we end up with a long horizon problem. That is, we plan for long-term gains. Assume we were training a Reinforcement

<sup>2</sup>Sutton & Barto, Sect. 8.4

<sup>3</sup>The Dyna-Q algorithm is one example of Dyna, see Lecture 7, Slides 7ff

Learning agent for a long-horizon problem. Could you think of a reason for which a method using short-horizon targets (cut-off at some horizon  $h$ ) might outperform a method using long-term horizon targets on this problem? To aid in answering, consider searching online for “planning horizon reinforcement learning model accuracy” and look for related work.

### **Part c**

What other parameters are relevant in Reinforcement Learning, and what considerations would help to determine their values?