

# Reinforcement Learning: An Introduction

## Solutions to the Second Edition

Joseph Early, University of Southampton

June 15, 2018

## 1 Introduction

### Exercise 1.1: Self-Play

The main consideration when the reinforcement learning (RL) algorithm plays itself is that its opponent is no longer constant - the opponent is also learning so its policy will change over time. As stated in the example, if the step size remains constant then the algorithm can remain competitive against an opponent that changes strategy over time, so a constant step size would be required in this instance. The overall level of both agents is likely to be greater than against a static opponent as both will learn from their strengths but also have their weaknesses exploited by their opponent. The final policies are likely to be better against any static opponent as a more general policy will have been found rather than one that just exploits a single opponent's weaknesses.

### Exercise 1.2: Symmetries

By using the symmetrical properties of tic-tac-toe, the number of states in the value function can be significantly reduced, which will speed up the learning process. However, this assumes that the opponent makes the same move according to the symmetry of the position, otherwise symmetrical states would have different values, in which case there is a potential weakness in the opponents policy that is not exploited. Since the state space for tic-tac-toe is relatively small, it is more important to look for an opponent's weakness by treating symmetric states as independent than using symmetries to reduce the number of states, especially when a draw is considered as bad as a loss.

### Exercise 1.3: Greedy Play

An RL agent that acted greedily would still be able to beat a constant player, however it would likely take longer to find a good strategy, and its policy may not be optimal, as once it had found a strategy that worked it would always exploit it rather than exploring for potentially better strategies. A greedy agent is likely to be more brittle as well - it would not have seen as many states as a non-greedy agent, therefore if it came up against a different opponent it would have less understanding of certain positions compared to a non-greedy agent that had explored those positions already.

## Exercise 1.4: Learning from Exploration

When learning from just greedy moves, the probabilities converge to the probability of winning from a state given that the best move is made. On the other hand, if learning occurs from all moves (greedy and exploratory), the probabilities converge to the probability of winning from a state given that the best move is made with probability  $1 - p$  and a random move is made with probability  $p$ , where  $p$  is the rate of exploration. If an agent continues to make exploratory moves after convergence, it is better to use the second probability set (the set produced by learning from both greedy and exploratory moves), as it is a more accurate representation of the agent's strategy so should result in more wins.

However, it is unnecessary for the final agent to keep making exploratory moves as it has already determined which moves are the best (assuming a constant opponent) - it is actually detrimental to keep exploring. An agent that doesn't explore once the probabilities have converged will achieve more wins than an agent that continues to explore, however the greedy agent should use the probability set produced by only learning from greedy moves as this represents the strategy that the agent actually follows. Therefore the best way to maximise wins for the final agent is to only learn from greedy moves and to reduce the exploration rate over time.

## Exercise 1.5: Other Improvements

One way to produce an agent that is a better overall tic-tac-toe player is to expose it to multiple opponents with different strategies. This means the agent will learn a more general strategy that will be better against any opponent rather than just the one it has played against, and that the agent will experience a greater number of states than it would do by playing against a single opponent. In this manner it could also learn strategies for playing as both noughts and crosses. This could be implemented by choosing an opponent at random (from a pool of possible opponents) for each game the agent plays.

Since the number of states is relatively small for tic-tac-toe, an alternative method would be to present each possible state to the opponent (assuming it plays deterministically) and establish a model for which action it will take in each state. It is then trivial to find an optimal policy to use against this opponent by adapting minimax search for the opponent's policy.

## 2 Multi-armed Bandits

### Exercise 2.1: Two-armed test bed

The probability of choosing the greedy arm is the probability of exploiting plus the probability of randomly choosing the greedy arm when exploring:

$$\begin{aligned} P(\textit{Greedy}) &= P(\textit{Exploit}) + P(\textit{Greedy}|\textit{Explore}) \\ &= 1 - \epsilon + \frac{\epsilon}{\textit{numArms}} \\ &= 0.5 + \frac{0.5}{2} \\ &= 0.75 \end{aligned} \tag{1}$$

### Exercise 2.2: Bandit example

It is impossible to determine if an action was definitely greedy from the information given - the agent could explore but happen to choose the greedy option at random, resulting in an action that appears greedy but was in fact exploratory. It is only viable to say if the action was exploratory or possibly exploratory, however one could argue that it is arbitrary whether the agent took the greedy action or explored and randomly selected the greedy action since the outcome is the same. The clearest way to evaluate which actions were definitely exploratory and which were not is to examine the action-value estimates over time:

Time	Arm 1	Arm 2	Arm 3	Arm 4	Greedy	Action	Reward
1	0	0	0	0	Any	1 - Possibly exploratory	1
2	1	0	0	0	1	2 - Exploratory	1
3	1	1	0	0	1 or 2	2 - Possibly exploratory	2
4	1	1.5	0	0	2	2 - Possibly exploratory	2
5	1	1.666...	0	0	2	3 - Exploratory	0

### Exercise 2.3: Greedy vs $\epsilon$ -Greedy

The  $\epsilon$ -greedy methods are clearly better than the greedy method, therefore it is only necessary to consider the two  $\epsilon$ -greedy methods. As the action-values converge, the optimal value selection also converges. When the action-values converge, the behaviour for the agent is given by Equation 1, which is the upper bound for the probability of selecting the greedy action. By substituting in the relative values of  $\epsilon$  for the two methods, it is found that the  $\epsilon = 0.1$  method is limited to 91% optimal action selection and the  $\epsilon = 0.01$  method is limited to 99.01% optimal action selection, meaning given infinite time steps the  $\epsilon = 0.01$  method will have a higher cumulative reward. Under the assumption that both methods spend a far longer time with action-values that have converged than they do converging, the improvement of the  $\epsilon = 0.01$  method over the  $\epsilon = 0.1$  method will tend towards  $\sim 1.088$ .

Add figure to show convergence.

In Section 2.5: Tracking a Non-stationary problem, it is mentioned that a constant step-size results in a weighted average of past rewards, and that the sum of the weights  $(1 - \alpha)^n + \sum_{i=1}^n \alpha(1 - \alpha)^{n-1} = 1$ , which the reader is encouraged to check. Below is a proof of the given statement:

$$\begin{aligned}
S &= (1 - \alpha)^n + \sum_{i=1}^n \alpha(1 - \alpha)^{n-1} \\
&= (1 - \alpha)^n + \alpha(1 - \alpha)^{n-1} + \alpha(1 - \alpha)^{n-2} + \cdots + \alpha(1 - \alpha) + \alpha \\
&= (1 - \alpha)(1 - \alpha)^{n-1} + \alpha(1 - \alpha)^{n-1} + \alpha(1 - \alpha)^{n-2} + \cdots + \alpha(1 - \alpha) + \alpha \\
&= (1 - \alpha)^{n-1}(1 - \alpha + \alpha) + \alpha(1 - \alpha)^{n-2} + \cdots + \alpha(1 - \alpha) + \alpha \\
&= (1 - \alpha)^{n-1} + \alpha(1 - \alpha)^{n-2} + \cdots + \alpha(1 - \alpha) + \alpha \\
&\vdots \\
&= 1 - \alpha + \alpha \\
&= 1
\end{aligned}$$

## Exercise 2.4: Arbitrary step-size parameter weighting

To determine the weighting of each prior reward for  $Q_{n+1}$  when the step-size parameters,  $\alpha_{n+1}$ , are not constant, a derivation similar to the one for constant step size can be used:

$$\begin{aligned}
Q_{n+1} &= Q_n + \alpha_n(R_n - Q_n) \\
&= \alpha_n R_n + (1 - \alpha_n)Q_n \\
&= \alpha_n R_n + (1 - \alpha_n)\alpha_{n-1}R_{n-1} + (1 - \alpha_n)(1 - \alpha_{n-1})Q_{n-1} \\
&= \alpha_n R_n \\
&\quad + (1 - \alpha_n)\alpha_{n-1}R_{n-1} \\
&\quad + (1 - \alpha_n)(1 - \alpha_{n-1})\alpha_{n-2}R_{n-2} \\
&\quad + (1 - \alpha_n)(1 - \alpha_{n-1})(1 - \alpha_{n-2})\alpha_{n-3}R_{n-3} \\
&\quad \vdots \\
&\quad + (1 - \alpha_n)(1 - \alpha_{n-1})\cdots(1 - \alpha_3)\alpha_2R_2 \\
&\quad + (1 - \alpha_n)(1 - \alpha_{n-1})\cdots(1 - \alpha_2)\alpha_1R_1 \\
&\quad + (1 - \alpha_n)(1 - \alpha_{n-1})\cdots(1 - \alpha_1)Q_1 \\
&= \alpha_n R_n + \sum_{i=2}^n \left( \alpha_{i-1} R_{i-1} \prod_{j=i}^n (1 - \alpha_j) \right) + Q_1 \prod_{i=1}^n (1 - \alpha_i)
\end{aligned}$$