# Self-Efficacy Update in Reinforcement Learning: Impact on Goal Selection for Q-learning Agents

Anonymous Author(s) Affiliation Address email

## Abstract

We introduce a dynamic self-efficacy learning rule and examine its impact on multi-1 2 goal selection in a grid-world. We model the Q-learning agent's self-efficacy as the integral of reward prediction errors (RPEs), allowing it to modulate the agent's 3 expectation of achieving the best possible future outcome. Initial simulation results 4 suggest that faster self-efficacy updates lead to higher overall reward accumulation, 5 but with increased variability in reaching the optimal goal. These findings indicate 6 that an optimal self-efficacy update rate, which can be learned through experience, 7 may strike a balance between maximizing performance and maintaining stability. 8

#### 9 1 Introduction

Multi-goal selection in reinforcement learning is challenging due to the need to balance exploration and exploitation across multiple objectives [Ecoffet et al., 2021]. Intrinsic motivation has been proposed as a mechanism to facilitate learning by integrating signals related to learning progress and competence [Barto, 2013, Colas et al., 2018]. Here, we model self-efficacy, the belief in one's ability to achieve desired outcomes [Bandura, 1997], in Q-learning agents and explore how varying self-efficacy updates impact goal selection and performance in a multi-goal grid-world environment.

### 16 2 Model

17 We model the *dynamics of self-efficacy* in response to feedback from the external environment and

18 propose that self-efficacy is a dynamic attribute, continuously shaped by action outcomes. The reward

<sup>19</sup> prediction error (RPE) is computed from a mixture of weighted best and worst future outcomes:

$$\delta_t = R_{t+1} + \gamma \cdot (w_t \cdot \max Q(s_{t+1}, a_{t+1}) + (1 - w_t) \cdot \min Q(s_{t+1}, a_{t+1})) - Q(s_t, a_t)$$
(1)

The *self-efficacy belief*,  $w_t$ , scales the highest possible future expected reward contingent on action, reflecting the agent's belief that it can successfully select the best action in the immediate future (as opposed to the worst one) [Gaskett, 2003, Zorowitz et al., 2020]. RPEs serve as the critical signal for updating both the action values and the self-efficacy parameter [Li and Radulescu, 2024]:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \cdot \delta_t \tag{2}$$

$$w_{t+1} = w_t + w_{LR+} \cdot \delta_t \tag{3}$$

#### 24 **3** Results and Discussion

In the multi-goal grid-world environment we tested, there are two target locations with rewards: Goal 25 1, located at (9, 9) with a reward of 4, and Goal 2, located at (0, 9) with a reward of 3. We evaluated 26 the impact of varying self-efficacy updates on the agent's goal selection behavior across 100 runs, 27 each run consisting of 200 episodes. With a low self-efficacy update rate, the agent reached Goal 1 in 28 51.22% of runs, accumulating an average reward of 636 with a standard deviation of 23.68. For the 29 medium self-efficacy update rate, The agent reached Goal 1 in 57.49% of runs and accumulated an 30 average reward of 660 with a standard deviation of 27.11. With the fast self-efficacy update rate, the 31 agent reached Goal 1 in 56.20% of runs and accumulated an average reward of 661 with a standard 32 deviation of 32.59. Overall, the medium self-efficacy update rate provided the best balance between 33 performance and stability, while fast updates led to increased variability in goal-selection behavior. 34



Figure 1: Impact of self-efficacy update rates on self-efficacy dynamics and goal selection behavior.

## 35 **References**

Albert Bandura. *Self efficacy: the exercise of control.* W. H. Freeman, New York (N.Y.), 1997. ISBN 978-0-7167-2850-4.

Andrew G. Barto. Intrinsic Motivation and Reinforcement Learning. In Gianluca Baldassarre
 and Marco Mirolli, editors, *Intrinsically Motivated Learning in Natural and Artificial Systems*,
 pages 17–47. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013. ISBN 978-3-642-32374-4

41 978-3-642-32375-1. doi: 10.1007/978-3-642-32375-1\_2. URL http://link.springer.com/

- 42 10.1007/978-3-642-32375-1\_2.
- 43 Cédric Colas, Pierre Fournier, Olivier Sigaud, Mohamed Chetouani, and Pierre-Yves Oudeyer.
  44 CURIOUS: Intrinsically Motivated Modular Multi-Goal Reinforcement Learning. *Proceedings of*
- the 36th International Conference on Machine Learning 2019, 2018. doi: 10.48550/ARXIV.1810.

46 06284. URL https://arxiv.org/abs/1810.06284. Publisher: arXiv Version Number: 4.

- 47 Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O. Stanley, and Jeff Clune. First return, then 48 explore. *Nature*, 590(7847):580–586, February 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/
- 49 s41586-020-03157-9. URL https://www.nature.com/articles/s41586-020-03157-9.

<sup>50</sup> Chris Gaskett. Reinforcement learning under circumstances beyond its control. *In Proceedings of the international conference on computational intelligence, robotics and autonomous systems*, 2003.

Jing Li and Angela Radulescu. Dynamic self-efficacy as a computational mechanism of mania emergence. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 46, 2024.

Samuel Zorowitz, Ida Momennejad, and Nathaniel D. Daw. Anxiety, Avoidance, and Sequential
 Evaluation. *Computational Psychiatry*, 4(0):1, 2020. ISSN 2379-6227.