

---

# Optimistic initialization of parameterized value functions using a neural network

---

## Abstract

Optimistic initialization of value functions is a popular approach to exploration in tabular reinforcement learning. However, it is rarely analyzed in deep reinforcement learning. We explore this problem through a parameterized value function using linear neural networks and compare our results to an existing popular learning algorithm.

## 1. Introduction

The persistent challenge of balancing exploration and exploitation is a fundamental challenge encountered by each reinforcement learning algorithm. Since this challenge arises everywhere and impacts the overall performance of algorithms, various methods to balance exploration or exploitation have been proposed. Among those approaches, one of the most fundamental and flexible is an optimistic initialization of the value functions. By simply setting the initial values greater than the reward maxima, one can force the agents to explore every state(-action) pair at least once at the early stage. The impact of optimistic initialization is noteworthy, particularly in scenarios involving tabular data, despite its straightforward nature. Although this simple technique is proven to provide tremendous advantages in a tabular context, its effect on deep nonlinear function approximation is yet to be discovered. This project aims to provide 1. a basic method to optimistically initialize the deep network; 2. a simple fix on the implementation so that optimistic values remain effective after a few gradient steps; and 3. empirical analysis.

## 2. Background

## 3. Research Question

This study aims to address the following three-part question: Can a parameterized value function using a neural network be optimistically initialized using reward signal shift and

normalization? How will the removal of the bias term from the neural network affect optimistic initialization? In a value estimation experiment, how will this initialization compare to its tabular counterpart using SARSA and semi-gradient SARSA with the same reward signal shift?

## 4. Experimental Design

We will be considering the minigrid library to run our experiments in three stationary environments. For the purposes of this paper, we will be focusing on the stationary environments Crossingenv, Distshiftenv and Lavagapenv.

The value estimation algorithm in our tabular setting is SARSA and for the parameterized setting we will be using semi-gradient SARSA, using an  $\epsilon$ -greedy behavior policy. The minigrid library API provides a discrete action space, along with encoding states in an image format for the parameterized setting and a coordinate format for the tabular setting.

We will perform SARSA using an epsilon-greedy policy to estimate the state-action value functions using the agent coordinates while initializing the state-action values optimistically. We will extend this idea to the parameterized case by adjusting the weights of the neural network so that it outputs an optimistic value for all state-action values. We will repeat the parameterized setting experiment by removing the bias term from the neural network and re-adjusting and performing the experiments accordingly. Comparing the results should give us better insight into the process of initializing parameterized value functions, its effect on expected return and the significance of the bias term in performance.

Furthermore, the mentioned environments were selected to due having simpler environment dynamics and a smaller state-space, making our experiment computationally less expensive.

## 5. Contributions

**Alireza Azimi:**

**Haruto Tanaka:**

**Henry Du:**

**Mashfique Zaman:**

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

**6. References**

To be added later.

**References**