

REINFORCEMENT LEARNING IN POKÉMON RED TO EXPLORE COMPLEX MULTI-REWARD ENVIRONMENTS

by

LIAM O'DRISCOLL
URN: 6640106

OCTOBER 30, 2023

Department of Computer Science
University of Surrey
Guildford, Surrey
England, United Kingdom
GU2 7XH

Project Supervisor: Sotiris Moschoyiannis

I declare that this dissertation is my own work and that the work of others is acknowledged and indicated by explicit references.

Liam O'Driscoll
October 30, 2023

© Copyright Liam O'Driscoll, October 30, 2023

0.1 Definitions

- ♣ Agent: The decision making mechanism receiving state information and performing chosen actions.
- ♣ Environment: The world in which the agent interacts with.
- ♣ State: A representation of the environment at the current timestep.
- ♣ Timestep: A value that increments after each action has passed since the start of the episode.
- ♣ Episode: An instance of the environment that the agent is interacting with.
- ♣ Action: The choice made by the agent in response to the state.
- ♣ Reward: The return value when an action is applied to a state.
- ♣ Reward Function: The mechanism in the environment that indicates how well the selected action is to achieving the goal of the environment.
- ♣ Policy: The decision making mechanism within the agent that decides the best action to perform given the state.

1 Introduction

The aim of this project is to develop a reinforcement learning agent to play Pokémon Red to compare the effectiveness of different styles of RL algorithms and their effectiveness. RL is an area of machine learning (ML) where agents make decisions and perform actions on states to achieve a goal.

1.1 Aims

- ✕ The aim of this project is to develop a RL agent to play Pokémon Red to compare the effectiveness of different styles of RL algorithms and their effectiveness to learn complex reward functions.

1.2 Objectives

- ◆ Research applications of RL to Pokemon and conduct a literature review on them
- ◆ Implement Pokémon Red game to be a suitable for training of different RL algorithms.
- ◆ Evaluate the performance of different RL algorithms used to train agents within the environment.
- ◆ Evaluate performance of agents to different forms of rewards functions.
- ◆ Recommend further developments to the project and applications to real world projects.

2 Literature Review

In RL, the decision making agent learns through experiences and 'trial and error'. Initially, it has a lack in understanding of the environment. However, through random action selection and the reward that it receives, it is able to learn an understanding of the environment. The agent is incentivized to maximise its reward and will aim to find actions that will yield more reward. This constant state, action and reward loop is what helps the agent improve by altering the policy after every cycle.

RL is different to common traditional machine learning techniques as it learns "how to map situations to actions-so as to maximize a numerical reward signal [1]." Agent's requirement to learn through experience and actions performed on current states not only affect the present, but also affect future states and actions are two characteristics which distinguishes itself from other forms of ML. It is also what makes Pokémon Red a suitable environment to apply this style of ML to.

RL has only been applied to every Atari game, where it surpassed the human benchmark for every game. However, these games lack long term randomness and present actions influencing future states. The game 'Pokémon Red' is an RPG game filled with various puzzles, non-linear world and a large amount of variance making each play through of the game unique while also require achieving the same goals. In addition, the game has 2 states the players is constantly in, the player is either in an overworld where they control their movement on a map or they are in a battle with another individual, where they control the actions of their monster.

Other similar works include is by Kalose et al [2] which applies RL to find the optimal battling strategy. Their work focuses on one aspect of the game and does not have a large enough search space to effectively apply rl techniques. Another similar work by Flaherty, Jimenez and Abbasi [3] applies RL algorithms A2C and DQN to play Pokémon Red. However, this piece of work does not go into enough detail about the comparison of different rl techniques to find the best method to train an agent to complete large complex environments with a large search space. I aim to extend their research in applying rl to Pokémon Red by applying more algorithms and various techniques rl that I will go into more detail in section 3.

I chose to apply rl to this specific environment, Pokémon Red, because of the benefits it holds with training and applications of research. This version of the game has the ability to speed-up the environment which allows for more timesteps to be completed and for the agent to experience more states. Another reason I chose to apply rl to this environment is because its complexity. The end goal of Pokémon Red is to defeat all the gym leaders and become the champion. However, to reach this goal the player must complete a series of smaller tasks which are not explicitly specified in the reward function. An example of this would be navigation in a 2-dimensional plane and solving puzzles along the way. The fact that

the agent is required to learn complete these tasks while completing the main goal can be applied and extended to real world applications.

3 Technical Overview

- Hyperparameter tuning to find optimal performance per experiment.
- Comparison of Gradient Descent and Value based models.
 - Value based:
 - * Proximal Policy Optimization
 - Gradient Descent:
 - * Actor-Critic Methods: A2C
 - * Deep Deterministic Policy Gradient
- Evaluating change in Q values to learning the optimal model using DQN
- Explore the benefit of applying meta learning.

4 Workplan

Month	Goals
October	<ul style="list-style-type: none"> • Rough structure of the report has been made. • Papers surrounding the project have been read (e.g., similar projects, algorithms that will be explored and technologies to be implemented) • Coding for the project is at its early stages. • Project Synopsis completed and submitted.
November	<ul style="list-style-type: none"> • Research and test which algorithms are applicable for comparison and applicable to project. • Draft introduction completed with a basic explanation of RL and how it is suitable for my environment. • Implementation of the Environment is complete
December	<ul style="list-style-type: none"> • Minimum viable product of code is achieved • Alter reward functions to give different incentives • Problem Analysis has been written • Design documentation and choice has been started
January	<ul style="list-style-type: none"> • Hyperparameter train sets of agents per algorithm • Train agents on different algorithms • Complete Design choice • Start evaluation of agents
February	<ul style="list-style-type: none"> • Any necessary extra agent training to be completed • First version of Report is at a Submittable state
March	<ul style="list-style-type: none"> • Debugging time for any potential issues • Review of draft report submission
April	<ul style="list-style-type: none"> • Consider completing Extension Objectives • Final report completed • Time allocated for debugging or potential issues
May	<ul style="list-style-type: none"> • Last final checks on final version of report • Time allocated for debugging or potential issues

5 References

- [1] R. S. Sutton, F. Bach, and A. G. Barto, *Reinforcement learning: An introduction*, Second. MIT Press Ltd, 2018.
- [2] A. Kalose, K. Kaya, and A. Kim, “Optimal battle strategy in pokemon using reinforcement learning,” *Web: <https://web.stanford.edu/class/aa228/reports/2018/final151.pdf>*, 2018.
- [3] J. Flaherty, A. Jimenez, B. Abbasi, B. Abbasi, J. Flaherty, and A. Jimenez, “Playing pokemon red with reinforcement learning,” 2021.