CSE574 Introduction to Machine Learning

Project 4

Bonus Part

Implementation of Two Environments from OpenAI's Gym Library

Project Report

Submitted By

Karan Manchandia

**Abstract**

In this bonus part, the task is to implement DQN using two environments from OpenAI's Gym Library. Gym is an open source interface to reinforcement learning tasks. It is a collection of test environments that we can use to workout reinforcement learning algorithms. The first environment I have chosen is CartPole-v1. The task is to balance a pole on a frictionless moving cart on which some force is applied.

The second environment used is a Atari environment from OpenAI's Gym Library called AlienNoFrameskip-v4. The task is to maximize the score in an alien game. In this environment, the observation is a screen RGB image which is represented as an array of shape (210,160,3). The action that is sent to the environment is applied and the state is returned without skipping any frames. The game resembles through the Arcade Learning Environment which uses Stella Atari emulator.

**Introduction:**

- **What I have Implemented?**

  I explore the application of DQN in 2 environments: the basic gym Cartpole environment where the goal is to balance a pole on a cart and the Atari Alien environment where the goal is to avoid the aliens and collect as much food particles as possible.

- **How the agent is trained?**

  - In the Cartpole environment, I use a Multilayer Perceptron (Neural Network) to estimate the Q-value function.

  - In the Atari environment, I have used a Convolutional Neural Network to estimate the Q–value function as a CNN is more suited for the Atari environment pixels.

- For the implementation, I have used stable baselines which is a reinforcement learning library that tries to match scikit-learn's easy to use interface. It contains pre-implemented agents including DQN. It is even possible to train a model using just one line of code.

- **How much time the does agent take to train itself ?**
  - For the bonus task 1 CartPole-v1 the DQN agent trains fairly quickly. It takes about 2 minutes to train.
  - For the bonus task 2 AlienNoFrameskip-v4, the agent takes over 2 hours to train.
  - This difference in time is mainly because of difference in complexity in classic control tasks and the Atari environment.

**Observations, Results, Screenshots and Graphs:**

**Bonus Part 1:**

**Screenshots:**



Figure 1: CartPole ffmpeg output

Figure 2: Bonus Part 1 Screenshot



Figure 3: Screenshot Showing Outputs

**Observation Table:**

| Timesteps | Best Mean Reward |
|-----------|------------------|
| 1954 | 22.36 |
| 4810 | 34.11 |
| 8871 | 69.12 |
| 12889 | 103.60 |
| 14950 | 120.51 |
| 18899 | 141.58 |

Graphs:



**Explanation:**

We can clearly see that as number of timesteps increases the rewards also increases. As the number of timsteps increases the agent moves form exploration to exploitation. As the agent began to exploit, it starts using q-values and has better moves, hence he reward also increases.
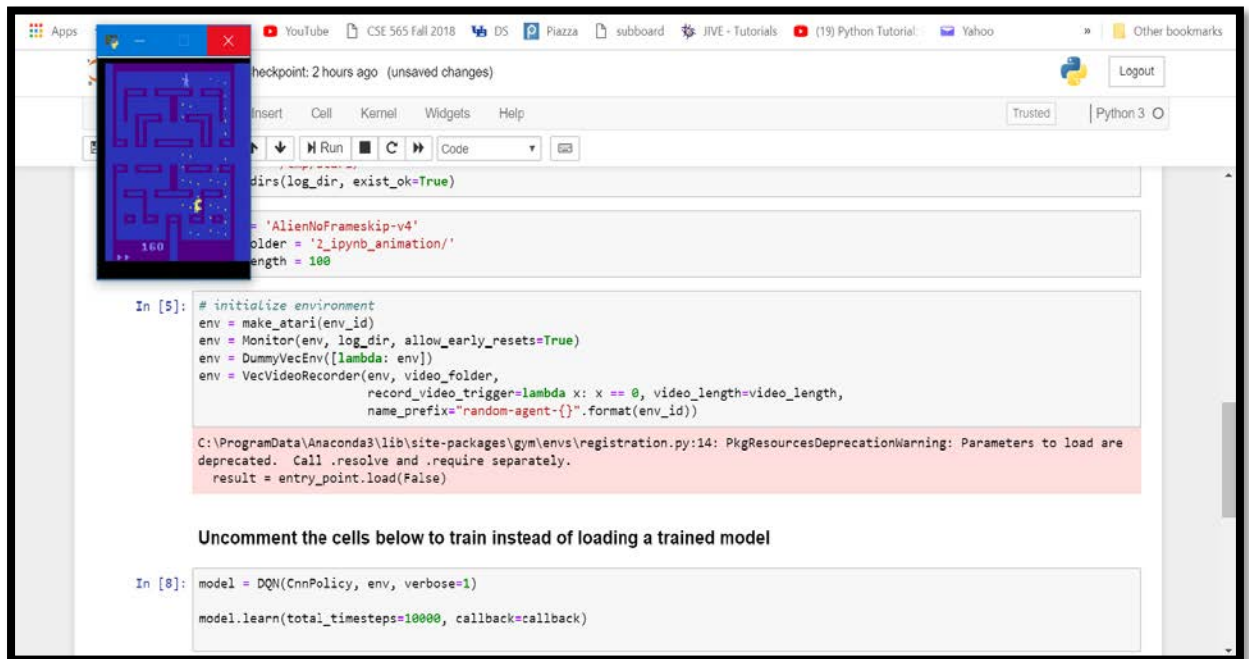
**Bonus Part 2 (AlienNoFrameskip-v4):**

**Screenshots:**
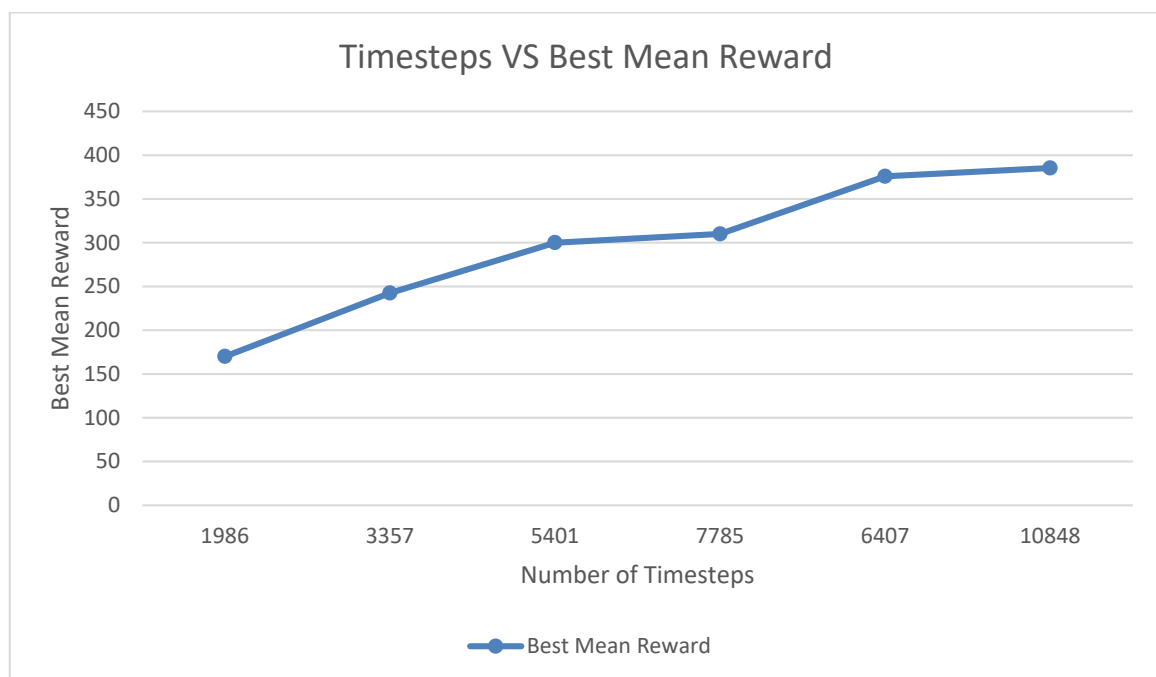


Figure 4 Bonus Task 2 Output Screenshot



Figure 5 Bonus Task 2 Output Screenshot

**Observation Table:**

| Timesteps | Best Mean Reward |
|---|---|
| 1986 | 170 |
| 3357 | 242.50 |
| 5401 | 300 |
| 7785 | 310 |
| 6407 | 375.83 |
| 10848 | 385.38 |

**Graphs:**



**Explanation:**

As seen above in Bonus Task 1, In task 2 the similar trend is seen. We can clearly see that as number of timesteps increases the rewards also increases. As the number of timesteps increases the agent moves form exploration to exploitation. As the agent began to exploit, it starts using q-values and has better moves, hence the reward also increases.

**References:**

- https://gym.openai.com/docs/
- https://mc.ai/openai-gym-environment-full-list/
- https://www.kdnuggets.com/2018/03/5-things-reinforcement-learning.html

- https://www.slideshare.net/AmazonWebServices/cmp305-deep-learning-on-aws-made-easycmp305

- https://medium.com/@m.alzantot/deep-reinforcement-learning-demysitifed-episode-2-policy-iteration-value-iteration-and-q-978f9e89ddaa