

# Deep reinforcement learning

Delfina Iriarte (1231682)

<sup>1</sup>University of Padova, Physics of Data

**Abstract.** *In this work, deep reinforcement learning techniques were applied using the Cartpole and the LunarLanding-v1 environment in PYTORCH. In particular, the model implemented was Q-learning which is a model-free reinforcement learning that has the goal to find an optimal policy, in the sense of maximizing the expected value of the total reward. Furthermore, the concepts of replay memory and target network were used.*

## 1. Introduction

Reinforcement learning tasks is about training an agent which interacts with its environment. The agent arrives at different scenarios known as states by performing actions leading to rewards. The goal of the agent is to interact with the emulator by selecting actions in a way that maximises future rewards. Q-Learning is a method of finding the optimal policies: the agent learns to associate a value to each **state-action pair**, which takes into account both the immediate reward for that action and the future rewards.

Deep reinforcement learning consists of algorithms that were based on function approximation i.e there was some function that approximates the Q values for actions. The Q-values founded for the next state can be used as an estimate of future values. In particular, the optimal action-value function obeys an important identity known as the Bellman equation. In practice, this basic approach was totally impractical and instead the update step can be consider as:

$$Q(s_t, a_t) \rightarrow Q(s_t, a_t) + \lambda \left( r_t + \gamma \max_a Q(s_{t+1}, a_t) - Q(s_t, a_t) \right) \quad (1)$$

where  $s_t, a_t$  were the state and the action at time  $t$ ,  $\lambda$  was the learning rate and  $\gamma$  the discount factor that takes into account how much future rewards matter.

Given the Q-values a policy is then needed. Many action selection techniques can be used for exploration such as epsilon greedy or softmax. Notice that exploration is crucial in the initial phases so that the agent find out as much as it can from the environment. If the agent was too greedy during the first episodes, it might get stuck in a local minima.

In this work, examples of reinforcement learning were implemented in PYTORCH. In particular, the environment studied were the Cartpole-v1 and the LunarLanding-v1 from the OPENAI GYM framework.

## 2. Methods

In this section, the environments, with its correspondent set of states, actions and rewards, are introduced. Furthermore, the network architecture consider for each of them are also presented.

### 2.1. Cartpole

The Cart-Pole environment consists of a cart that moves along the horizontal axis and a pole that was anchored on the cart. At every step the observations are: position, velocity, angle and angular velocity. The set of actions of the cart were only two: move to the left or move to the right. In this environment, the reward was given as long as the pole was balance in which the agent gets +1. An episode was over as soon as the pole falls beyond a certain angle or the cart strays too far off to the left or right. An example of the rendered environment was shown in Figure 7.

The agent acts in the environment and explore the world. All the experience of the agent were stored in the **replay memory** inside a buffer space and the model was train by sampling out the batches of experiences randomly or prioritized experiences. By sampling from it randomly, the transitions that build up a batch were decorrelated. It has been shown that this greatly stabilizes and improves the DQN training procedure. Furthermore, **target network** was used to compute  $Q(st + 1, at)$  for added stability. The target network has its weights kept frozen most of the time, but was updated with the policy network's weights every so often.

As an action selection, the **softmax with decreasing temperature** was used. Nevertheless, when the temperature was set to zero, the epsilon greedy exportation was applied which is a common choice in this models. In particular, an exponential decay for the exploration profile was used with an adjust parameter called `init_value`. The stochastic gradient descent (SGD) but without momentum was used as optimizer of the network. Finally, it should be also mentioned that the reward also has a linear penalty when the cart was far from the center since it was training to escape from the screen.

The network architecture of the model was shown in Figure 1. It takes as input the four state space dimension, then it pass through two fully connected layers. As activation function the `tanh` function has been used.

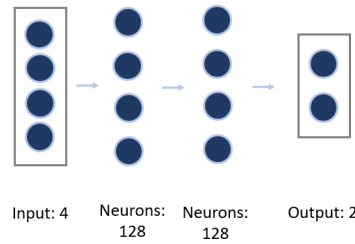


Figure 1: Model architecture for the cartpole environment

In order to speed up the convergence of the model, several exploration profiles were study with different starting temperature. Finally, the model performance was also compared to one obtained by a random agent.

## 2.2. Cartpole using screen pixels

Deep reinforcement learning (DRL) has been leading to state-of-the-art performance to learn control policies for a wide range of applications. The CartPole task (Section 2.1) was designed so that the inputs to the agent were 4 real values representing the environment state (position, velocity, etc.). However, neural networks can solve the task purely by looking at the scene. Strictly speaking, states were present as the difference between the current screen patch and the previous one. This will allow the agent to take the velocity of the pole into account from one image. Unfortunately this does slow down the training, because we have to render all the frames. Therefore, as **input the screen** it was provided. Notice that the image was quite large ( $800 \times 600 \times 3$ ), therefore the following manipulation were done:

1. Convert the image in grayscale considering only one channel.
2. Crop the image removing empty undesired screen.
3. Rescale up to  $1 \times 85 \times 200$

The network model was a convolutional neural network as shown in Figure 2. As activation function the ReLu has been used for the convolutional layer and tanh for the fully connected ones.

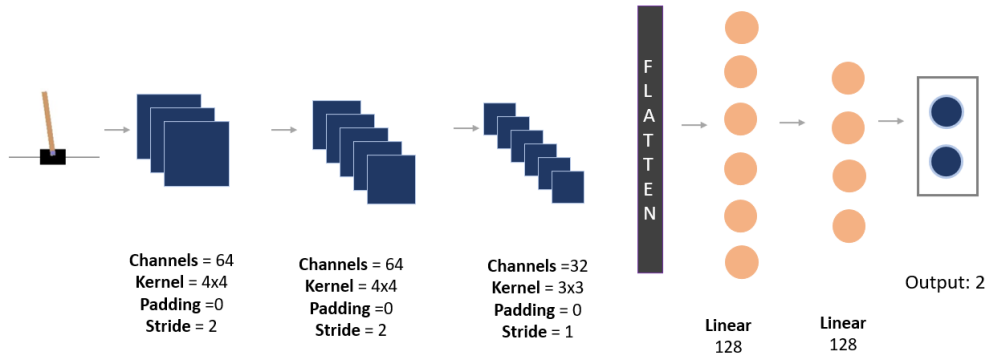


Figure 2: Model architecture for the cartpole environment using screen as input.

### 2.3. Landing A Rocket

In this section, Q-learning agent was tasked to learn the task of landing a spacecraft on the lunar surface given by the environment called `LunarLanding-v1`. A picture of the rendered environment was shown in Figure 8. The landing pad was always at coordinates (0,0). The states of the systems were the following:

1.  $s[0]$ : horizontal coordinate
2.  $s[1]$ : vertical coordinate
3.  $s[2]$ : horizontal speed
4.  $s[3]$ : vertical speed
5.  $s[4]$ : angle
6.  $s[5]$ : angular speed
7.  $s[6]$ : 1 if first leg has contact, else 0
8.  $s[7]$ : 1 if second leg has contact, else 0

whereas the the action available were the following discrete ones:

1. do nothing
2. fire left orientation engine
3. fire main engine
4. fire right orientation engine

Reward for moving from the top of the screen to the landing pad and zero speed was about 100-140 points. If the lander moves away from the landing pad it loses reward. The episode finishes if the lander crashes or comes to rest, receiving an additional -100 or +100 points. Each leg with ground contact was +10 points. Firing the main engine was -0.3 points each frame. Firing the side engine was -0.03 points each frame. Solved was 200 points.

Again, a random agent will be train in order to make comparison with the deep Q agent. The network architecture for this model was similar to the one of Section 2.1, but in this occasion was composed by 3 linear layers as shown in Figure 3. As activation function the tanh function was used. Additionally, a linear bonus for going downward was consider in order to avoid the space to fly indefinitely.

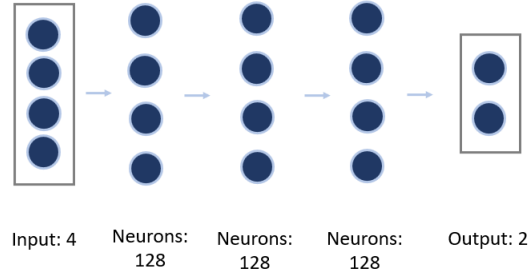


Figure 3: Model architecture for the lunar landing environment.

### 3. Results

In this section, the main results were present and discuss.

#### 3.1. Cartpole

Figure 4 presents the training score per iteration for a random agent and its frequency. As it can be observed, it's impossible to solve the environment using this approach. Even tough sometimes was lucky enough to get a reward of almost 75, it can be noticed that the agent was not learning from their experience. In fact, the average performance was of  $\sim 20$ .



Figure 4: Training score and frequency of the rewards for the Cartpole environment when training a random agent.

The training score and its frequency for several hyperparameters were shown in Figure 9 and Figure 10. It can be observed that the trend line was positive, since the performances increases over time and indeed the model was able to learn. The accuracy of the model highly depend on the hyperparameter choosen. In particular, the best parameter corresponds to the model which wins first which occurs in 436 number of episodes. The best parameters obtained were: `'gamma': 0.94`, `'lr': 0.062`, `'target_net_update_steps': 5`, `'batch_size': 128`, `'bad_state_penalty': 0` and `init_time: 25`. Finally, it should stress out that the training score oscillates a lot.

#### 3.2. Cartpole using screen pixels

Figure 5 shows the evolution of the reward of the cartpole. It can be observed that the model was not able to win the game, the average test score was around 102. In fact, it was resemble to a random agent in the sense that there was not any positive trend and the algorithm was not learning.

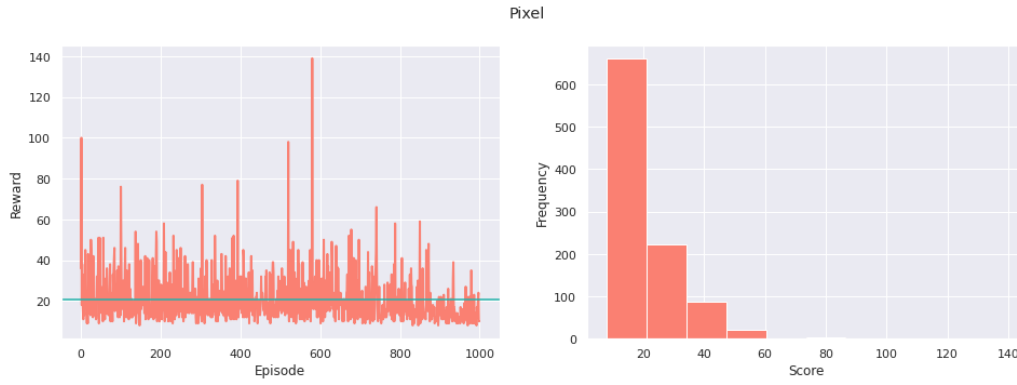


Figure 5: Training score and frequency of the rewards for the CartPole with screen input when training a random agent

### 3.3. Landing A Rocket

As in Section 3.1, the random agent for the model was not able to achieve high performance as shown in Figure 6, for the same reasoning as before. In fact, the rewards were on average negative.



Figure 6: Training score and frequency of the rewards for the lunar landing environment when training a random agent

Figure 11 presents the training score of the several parameters try in the model. As it can be observed, the environment was solved in almost all the cases. In particular for ' $\gamma$ ': 0.98, ' $lr$ ': 0.010, ' $target\_net\_update\_steps$ ': 10, ' $batch\_size$ ': 128, ' $bad\_state\_penalty$ ': 0 and  $init\_time$ : 9 it can be observed that the score with more frequency correspond to 200 which was the solved game.

## 4. Conclusion

In conclusion, deep reinforcement learning seems a challenging and promising method in order to solve "games" by learning from the environment. Despite the efforts, most of the games were successfully solved by tuning the hyper parameters. Unfortunately, one of the main drawback for the model was the RAM consuming, specially when dealing with pixel images.

---

# 5. Appendix

CartPole-v0 Environment



Figure 7: Rendered environment of the Cartpole

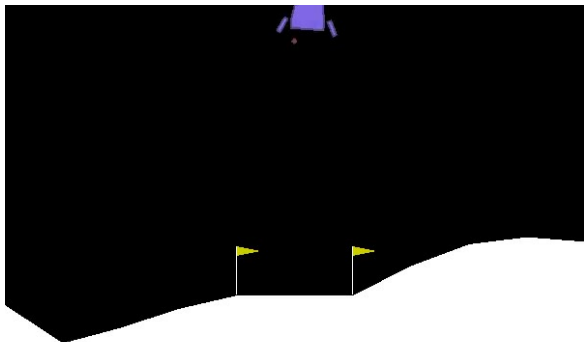


Figure 8: Rendered environment for the lunar landing environment.



Figure 9: Training score and frequency for the several hyperparameters for the Carpole environment



Figure 10: Training score and frequency for the several hyperparametes for the Carpole environment.



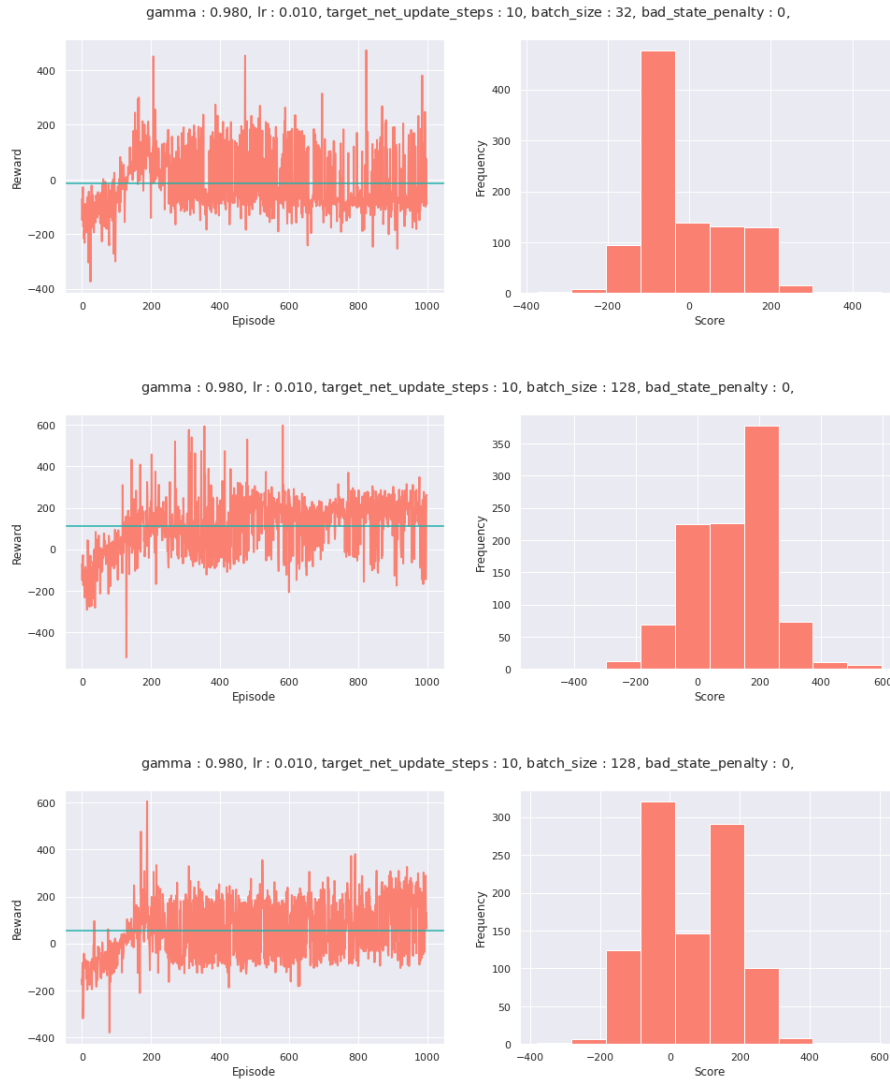


Figure 11: Training score and frequency for the several hyperparameters for the LunarLanding environment.