

Capstone Project Report: Learn to play Atari through Deep Reinforcement Learning

1 Definition

1.1 Project Overview

In reinforcement learning (RL) an agent interacts with an environment over a number of discrete time steps. At each time step the agent experiences a state s of the environment and decides to perform an action a . The action to be performed in a particular state is defined by the policy $\pi(a|s)$, and the quality of the policy is expressed in terms of the return R_t through the action-value function $Q^\pi(s, a) = \mathbb{E}[R_t | s_t = s, a]$. In Q-learning the learning task involves finding the optimal value function $Q^*(s, a)$, which leads to the maximum possible return. Once $Q^*(s, a)$ is known, the strategy is to act greedily with respect to a , or in other words to select at each step the action that maximise $Q^*(s, a)$. If the state space is very large a function approximator, typically a neural network (NN), is used to model the optimal policy $Q^*(s, a) \approx Q(s, a, \theta)$. This type of learning agent, which is called deep Q-network (DQN), led to the first algorithm that outperformed humans in playing Atari [1].

Nevertheless DQN can suffer from a problematic convergence and large training times. Convergence issues have been addressed through experience replay, and by adopting a separated target network to prevent training on correlated samples [2]. Still, this solution does not guarantee good performances in all the different RL problems. More recently policy gradient (PG) methods have been revived as they have shown for some applications superior convergence while reducing training time compared to DQN [3]. In this class of algorithms a NN is used to approximate the policy that maximises the expected reward $\pi^*(a|s) \approx \pi(s, a, \theta)$. Policy parameters learning is based on the gradient of the performance measure with respect to the policy parameters. PG methods however also do not work well universally and may require long time to find the correct hyperparameters. Perhaps the simplest advantage that policy parameterization may have over action-value is that the policy may be a simpler function to approximate compared to the action-value function. In fact, problems can have different complexity for their policies or action-value functions. For some, the action-value function is simpler and thus easier to approximate. For others, the policy is simpler. In the latter case a policy-based method will typically learn faster and yield a superior asymptotic policy [4].

1.2 Problem Statement

The aim of this project is to train a RL agent to play the Atari 2600 game Pong [5] through two different reinforcement learning approaches described in Section 1.1. In Pong an agent and his opponent control a paddle, and the game consists of bouncing a ball past the other player. At each time step the agent observes a pixel image and chooses between two actions, Up or Down, then he receives a reward based on whether or not the ball went past the opponent. The reward is +1 if the ball went past the opponent, -1 if it went past the agent. In the original DeepMind implementation [1, 3] the states consist of the last 4 84×84 pixels frames observed

by the agent, which results on $256^{84 \times 84 \times 4} \approx 1.4 \times 10^{67970}$ different combinations. In such a large state space the state-value function $Q^\pi(s, a)$ or the policy $\pi(s, a)$ cannot be stored in a table, but needs to be approximated by a NN. Such NN often requires several layers to beat the opponent consistently, which is the reason why this approach is called deep reinforcement learning.

1.3 Metrics

Each Pong episode is made of multiple consecutive games, and it finishes when one of the opponent reaches a score of 21. An agent who is not able to learn, or select actions at random, will lose the game most of the time. As a result, the return, i.e. the sum rewards at the end of an episode, will be always near to the minimum of -21. When the agent learns the action-value function, or the policy, it will start to outperform the opponent for an increasing number of games and the return will increase accordingly. Deep RL algorithms require in general very long training times to be able to beat consistently the agent, and are in general run on GPUs to speed up the calculations. The outstanding results in [3] were achieved after playing 10 millions episode, which translates into weeks of computational time. However, all the analyses in this project were run on a MacBook CPU 2.8 GHz 8 GB RAM. Since the scope of this project is to evaluate the performance of different algorithms, and not to outperform consistently the opponent, the training was limited to 2000 episodes. Therefore the performance of an algorithm is based on the evolution of the return along 2000 episodes played.

2 Analysis

2.1 Data Exploration and Visualisation

In a reinforcement learning problem the agent is not trained on a labelled data set, as in supervised learning, but instead he is required to understand the quality of his actions based on the rewards that he receives from the environment. In this project OpenAI Gym is used to replicate Atari 2600 Pong environment [5]. Each input is an RGB image as shown in Figure 1, which corresponds to an array of shape (210, 160, 3), and the environment returns a reward of +1 if the ball went past the opponents, -1 if the ball is missed, or 0 otherwise.

The agent can choose between 6 discrete actions labelled from 0 to 5, in which (0,1) are idle actions that don't have any affect, (2,4) make the paddle go up, and (3,5) make the paddle go down [?]. Once an action is selected it is repeatedly performed for a duration of k frames, where k is uniformly sampled from 2,3,4. This introduces some degree of randomness in the agent behaviour, as the same policy leads to slightly different agent behaviours when repeated in different episodes.

2.2 Algorithm and Techniques

2.2.1 DQN

The DQN algorithm has been introduced by the Google Deepmind paper in 2015 [2]. A DNN is used to approximate the optimal action-value function:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \gamma^{T-t} r_T | s, a, \pi] \quad (1)$$

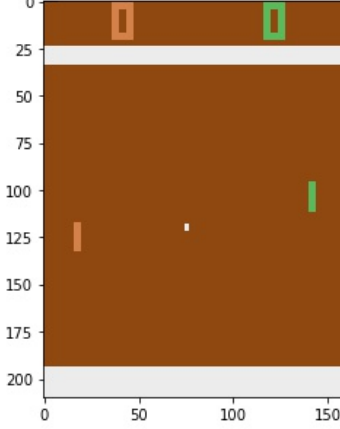


Figure 1: Pong frame.

which corresponds to the maximum expected return that can be achieved at any point in time following the policy π , under the hypothesis of discounted future rewards parametrised by $\gamma < 1$. The optimal action-value function can be found iteratively from the Bellman equation, which states that if $Q^*(s', a')$ was known at the next state s' for all actions a' , then the optimal greedy strategy is to select the action a' that maximizes the expected value of $r + Q^*(s', a')$:

$$Q^*(s, a) = r + \gamma \max_{a'} Q^*(s', a') \quad (2)$$

An action is selected according to an ϵ -greedy strategy in which the greedy policy is followed with $1-\epsilon$ probability, while a random action is selected with ϵ probability. A NN approximator of the optimal action value-function $Q^*(s, a) \approx Q(s, a, \theta)$ is trained on target values $y_i = r + \max_{a'} Q(s, a, \theta_{i-1})$ by defining a loss function:

$$L(\theta_i) = (y_i - Q(s, a, \theta_i))^2 \quad (3)$$

this means that the NN is updated using as target values its own outcomes obtained with current weights θ_{i+1} , which are held fixed during the training step. Although this approach is proven to converge as $t \rightarrow \infty$, in practice it is observed that correlation that between $Q(s, a, \theta_i)$ and the target value represents an issue. For this reason a strategy called *experience replay* is adopted in which the observed transitions (s_t, a_t, r_t, s_{t+1}) are stored into a memory D . The weight update is then performed at the end of each episode on a random batch of transitions sampled from D . A summary of the DQN algorithm is shown in Figure 2

2.2.2 Policy Gradient

Policy gradient refers to a class of algorithms that originate from 1992 Williams' REINFORCE algorithm [7]. Policy gradient method approximates the optimal policy $\pi^*(a|s) \approx \pi(s, a, \theta)$, instead of the action-value function $Q(s, a)$ used in DQN, and uses an unbiased estimate of the policy gradient to perform the weight update. Given the objective function $J(\theta) = V_{\pi_\theta}(s_0)$ that measures the quality of policy π_θ , the unbiased estimate of its gradient is given by [4]

$$\nabla J(\theta) = \mathbb{E}[R_t \nabla_\theta \log \pi(s, a, \theta)] \quad (4)$$

which is used to perform the weight update by gradient ascent, since now we want to maximise $J(\theta)$.

```

Initialise action-value function parametrization  $Q(s, a, \theta)$  with random weights  $\theta$ 
Initialise memory  $D$ 
for  $episode = 1, \dots, N$  do
  repeat
    Observe state  $s_t$  and perform action according to  $\varepsilon$ -greedy scheme
    Observe episode reward  $r_t$  and next state  $s_{t+1}$ 
    Store transition  $(s_t, a_t, r_t, s_{t+1})$  into  $D$ 
    Sample random mini-batch of transitions  $(s_k, a_k, r_k, s_{k+1})$  from  $D$ 
    if  $s_{t+1}$  is terminal then
       $y_k = r_k$ 
    else
       $y_k = r_k + \gamma \max_{a'} Q(s_{t+1}, a', \theta)$ 
    end if
    Perform gradient descent  $\theta \leftarrow \theta + \alpha(y_k - Q(s, a, \theta)) \nabla Q(s, a, \theta)$ 
  until episode is finished
end for

```

Figure 2: DQN algorithm

In this work the basic implementation by Andrej Karpathy [8] and described in Figure 3 will be used. A pre-defined number of episodes are played and stored in memory, and a discounted total reward is computed for each observed transition

$$r_{d,t} = \sum_{k=0}^{\infty} \gamma^k r_{t+k} \quad (5)$$

the discounted rewards are normalised, and weight update is performed on the batch of samples currently stored in memory

```

Initialise policy parametrization  $\pi(s, a, \theta)$  with random weights  $\theta$ 
Initialise batch memory  $D$ 
for  $episode = 1, \dots, N$  do
  repeat
    Observe state  $s_t$  and perform action  $a_t \sim \pi(s, a, \theta)$ 
    Observe reward  $r_t$  and compute discounted reward  $r_{d,t}$ 
    Store  $(s_t, a_t, r_{d,t})$  into  $D$ 
  until episode is finished
  if  $episode \bmod batch\ size = 0$  then
    Perform weight update  $\theta \leftarrow \theta + \alpha r_{d,t} \nabla \ln \pi(s_t, a_t, \theta)$ 
    Empty  $D$ 
  end if
end for

```

Figure 3: Policy gradient algorithm

2.3 Benchmark

The benchmark used to evaluate the performance of an algorithm is represented by the performance of an agent that selects actions uniformly at random. Figure 4 shows that the moving average of the return remains between -21 and -20, that is the opponent is consistently winning most the games in an episode and the performance of the agent is not improving. It is expected that a *learning* agent starts initially with similar performance to the random agent, but that he will improve with the number of games played.

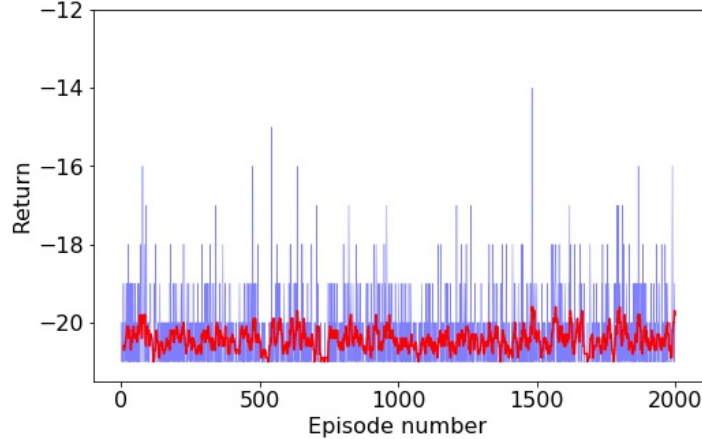


Figure 4: Performance of a random agent, in blue the episode return and in red the moving average on a window of 10 episodes.

3 Methodology

3.1 Data preprocessing

Figure 5a suggests that only the portion of the image between the two grey bands is needed to learn the game, and that colors do not contain any valuable informations. Therefore, a function is defined to pre-process the frames, which crops the central part of the frame and convert the image to a binary array with 0 assigned to the background and 1 to both the ball and the paddles. In addition, the image resolution is such that the geometries are preserved if the pixels are down-sampled by a factor of 2. The resulting frame which is shown in Figure 5b contains the same amount of information as the original one, but reduces the complexity of the NN used by the learning agent.

A learning agent needs information of the current state as well as the previous ones in order to improve his performance. This because the correct action can be chosen only if the trajectory of the ball is known or, in other words, if both velocity and position are available. In the DQN implementation of Google Deepmind [2] the authors have used a stack of the 4 previous images as current state. This solution leads to a computationally expensive model which was well suited for their CNN architecture. A simplified solution, which has been introduced in [8], and is widely adopted for testing codes, is to let the agent observe the difference between the current and the previous frame. The difference image which is shown in Figure 5c indicates that the

opponent (on the left) has moved downward, the player (on the right) has moved upward, and the ball is going towards the opponent.

In Section 2.1 it was observed that out of the 6 actions that the agent can perform in the Pong OpenAI version, only 2 are needed to play a game. This allows to simplify the output of the NN into a binary index (0,1) that maps to action label 2 to move the paddle up, and action label 3 to move the paddle down. Although this is not expected to reduce the computational complexity it allows to simplify the code.

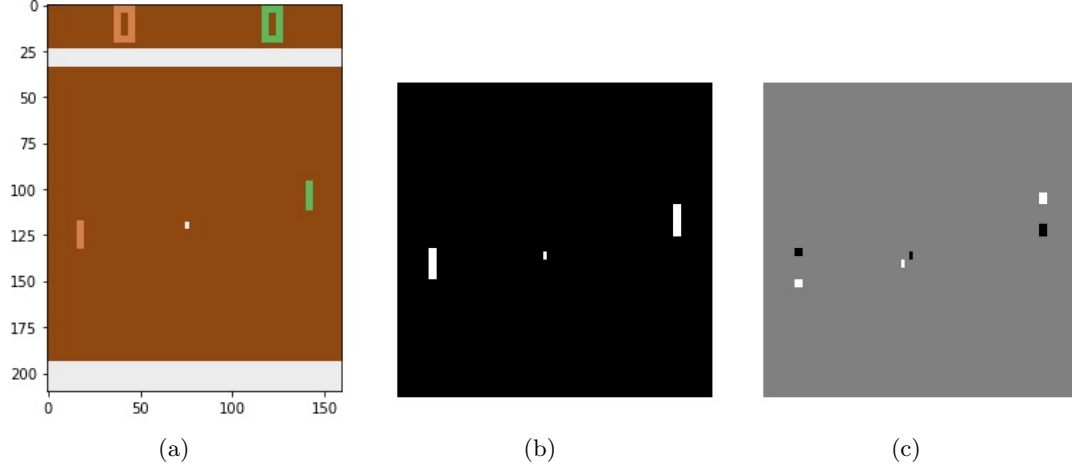


Figure 5: Frame preprocessing: (a) current original frame, (b) binary image after cropping and down-sampling, and (c) difference between next frame and current frame.

3.2 Implementation

Following a similar approach to [6] a RL algorithm is organised into 3 Python classes:

- the **Agent** who performs actions, observes the environment, and learns from the feedback of the environment;
- the agent's **Brain** which contains the definition of the NN architecture, training, prediction, and weight update procedures
- the **Environment** in which the agent learns, which returns the new state and the reward that follow an agent action in a loop that repeat for a large number of episodes

The initial NN architecture tested for both the DQN and the PG methods was the same as in [1] composed of 2 convolutional layers and a fully connected layer, and used as inputs stacks of 2 consecutive frames (instead of 4 frames as in the original paper). This configuration however required very large computational times which posed limitations to the testing and development of the codes. For this reason the same simplified network architecture in [8] was adopted for both the DQN and the Policy Gradient algorithm which consisted of a single fully connected hidden layer of 200 units. The algorithms which are described in detail in the next section were implemented using *Keras* and *Tensorflow* Python libraries.

3.2.1 DQN algorithm

The **Agent** class stores samples into the memory, update the ϵ parameter, maps the NN outputs to the action label, and implements the **replay** function which calculates the NN target using

Equation 2 and performs the network weight update. The target network is used to assign to the agent network the target Q-value (y_i in Equation 3) for the current action. A more efficient implementation which use a single network to perform the learning step [9], was also tested but this didn't reduce significantly the training time.

```

1 class Agent:
2     def __init__(self, stateCnt, actionCnt):
3         self.stateCnt = stateCnt
4         self.actionCnt = actionCnt
5         self.memory = deque(maxlen = REPLAY_MEMORY)
6         self.brain = Brain(stateCnt, actionCnt)
7
8     def action_index(self, s):
9         """ Select action index action according to epsilon-greedy scheme """
10        if random.random() < self.epsilon:
11            index = random.randint(0, self.actionCnt-1)
12        else:
13            index = np.argmax(self.brain.predictOne(s))
14        return index
15
16    def observe(self, sample):
17        """ Add sample to memory, decrease epsilon, update target network"""
18        self.memory.append(sample)
19        self.steps += 1
20        if self.epsilon > FINAL_EPSILON:
21            self.epsilon -= (INITIAL_EPSILON - FINAL_EPSILON) / EXPLORATION
22        # update target network
23        if self.steps % UPDATE_TARGET_FREQUENCY == 0:
24            self.brain.updateTargetModel()
25
26
27    def replay(self):
28        """ Sample experience from memory and update weights """
29        if len(self.memory) < BATCH_SIZE:
30            return
31        # sample from memory queue
32        batch = np.array(random.sample(self.memory, BATCH_SIZE))
33        states = np.vstack(batch[:,0])
34        actions_index = np.vstack(batch[:,1])
35        rewards = np.vstack(batch[:,2])
36        next_states = np.vstack(batch[:,3])
37        terminals = np.vstack(batch[:,4])
38
39        # current state predictions from AGENT network
40        Q = self.brain.predict(states)
41        # next_state predictions from TARGET network
42        next_Q = self.brain.predict(next_states, target=True)
43        # target Q values - if terminal: target_Q = rewards
44        target_Q = rewards + GAMMA * np.max(next_Q, axis=1) * np.invert(
terminals)
45        # assign target values to actions
46        Q[range(BATCH_SIZE), actions_index] = target_Q
47        # AGENT network update
48        self.brain.train(states, Q)
49        self.Q_sa = next_Q

```

Listing 1: DQN Agent definition

The Brain class defines the NN model of the action-value function $Q(s, a, \theta)$ implemented in *Keras*, which uses a RMSprop optimizer and a custom Huber loss function to clip the gradient [9]. The output of the network is a binary array (0,1) which is then used to select the paddle

actions Up and Down.

```
1 class Brain:
2     def __init__(self, stateCnt, actionCnt):
3         self.stateCnt = stateCnt
4         self.actionCnt = actionCnt
5         self.model = self._createModel()
6         self.model_ = self._createModel() # target model
7
8     def _createModel(self):
9         """ Implement NN to approximate action-value function  $Q(s,a)=Q(s,a,w)$  """
10
11         model = Sequential()
12         model.add(Dense(output_dim=200, activation='relu', input_dim=self.stateCnt))
13         model.add(Dense(output_dim=actionCnt, activation='linear'))
14
15         opt = RMSprop(lr=LEARNING_RATE, decay=DECAY_RATE)
16         model.compile(loss=huber_loss, optimizer=opt)
17         return model
18
19     def train(self, x, y, epochs=1):
20         """ Trains the model for a fixed number of epochs """
21         self.model.fit(x, y, batch_size=BATCH_SIZE, epochs=epochs, verbose=0)
22
23     def predict(self, s, target=False):
24         """ Batch prediction from agent network or target network """
25         if target:
26             return self.model_.predict(s)
27         else:
28             return self.model.predict(s)
29
30     def predictOne(self, s, target=False):
31         """ Predict one action from model """
32         return self.predict(s.reshape(1, self.stateCnt), target=target).flatten()
33
34     def updateTargetModel(self):
35         """assign weights to target Network """
36         self.model_.set_weights(self.model.get_weights())
```

Listing 2: DQN Brain definition

The `Environment` class implements the DQN algorithm shown in Figure 2 in which the agent observes a state, represented by the difference between the current and the previous frame, performs an action and observes the new state and reward. Then a transition is stored in memory and the network weight update is performed through the `replay` function at the end of each step.

```
1 class Environment():
2     def __init__(self, game):
3         self.game = game
4         self.env = gym.make(game)
5
6     def run(self, agent):
7         """ DQN algorithm """
8         obs = self.env.reset()
9         prev_image = preprocess(obs)
10        state = np.zeros(stateCnt)
11        R = 0
12        while True:
```



```

13         # agent performs action: labels:(0,1) -> actions:(2,3)
14         action_index = agent.action_index(state)
15         action = action_index+2
16         # environment return new state and reward
17         obs, reward, done, _ = self.env.step(action)
18         curr_image = preprocess(obs)
19         # agent learn from difference between current and previous frame
20         next_state = curr_image - prev_image
21         # store state in memory and update epsilon
22         agent.observe((state, action_index, reward, next_state, done))
23         # sample batch from memory and perform weight update
24         agent.replay()
25         # update state and frame
26         state = next_state
27         prev_image = curr_image
28         R += reward
29         if done:
30             break

```

Listing 3: DQN Environment definition

3.2.2 Policy Gradient Algorithm

The `Agent` class contains the memory of transitions, which is called by the `learn` function during weight update, and a function to clear the memory after each update on the memory batch. It implements also the function that calculate for each reward the corresponding normalized discounted reward until the end of the current game (and not until the end of a full episode).

```

1 class Agent:
2     def __init__(self, stateCnt):
3         self.stateCnt = stateCnt
4         self.memory = {'states': [], 'action_labels': [], 'rewards': []}
5         self.brain = Brain(stateCnt)
6
7     def clear_memory(self):
8         """Empty agent memory after an update"""
9         for key, _ in self.memory.items():
10             self.memory[key] = []
11
12     def action_label(self, state):
13         """Perform action"""
14         action_probs = self.brain.predict(state)
15         label = 1 if np.random.uniform() < action_probs[0,0] else 0
16         return label
17
18     def learn(self):
19         """Get samples of an episode from memory and update policy"""
20         states = np.vstack(self.memory['states'])
21         actions = np.vstack(self.memory['action_labels'])
22         rewards = np.vstack(self.memory['rewards'])
23         self.brain.update_model(states, actions, rewards)
24
25     def discount_rewards(self, r):
26         gamma = 0.99
27         """ take 1D float array of rewards and compute discounted reward """
28         discounted_r = np.zeros_like(r)
29         running_add = 0
30         # for t in reversed(range(0, r.size)):
31         for t in reversed(range(0, len(r))):

```

```

32         if r[t] != 0: running_add = 0 # reset the sum, since this was a
game boundary (pong specific!)
33         running_add = running_add * gamma + r[t]
34         discounted_r[t] = running_add
35         discounted_r -= np.mean(discounted_r)
36         discounted_r /= np.std(discounted_r)
37         return discounted_r.tolist()

```

Listing 4: Policy Gradient Agent definition

The `Brain` class contains the NN model of the policy $\pi(s, a, \theta)$ implemented in *Tensorflow*. The loss function used to calculate the gradient in Equation 4 is obtained by multiplying the discounted reward and the Tensorflow function `softmax_cross_entropy_with_logits` which returns $-\log(\pi(a_t|s_t))$. The NN returns the probability of an action `self.probs` through the sigmoid function. The same pre-process function used in the DQN algorithm is also included here.

```

1 class Brain:
2     def __init__(self, stateCnt):
3         self.stateCnt = stateCnt
4         self.states_ph = tf.placeholder(dtype=tf.float32, shape=(None, self.
stateCnt))
5         self.actions_ph = tf.placeholder(dtype=tf.float32, shape=(None,1))
6         self.rewards_ph = tf.placeholder(dtype=tf.float32, shape=(None,1))
7         self.model = self._createModel()
8         self.sess = tf.Session()
9         self.saver = tf.train.Saver()
10        self.sess.run(tf.global_variables_initializer())
11
12    def preprocess(self, I):
13        """ prepro 210x160x3 uint8 frame into 6400 (80x80) 1D float vector """
14        I = I[35:195] # crop
15        I = I[:,::2,::2,0] # downsample by factor of 2
16        I[I == 144] = 0 # erase background (background type 1)
17        I[I == 109] = 0 # erase background (background type 2)
18        I[I != 0] = 1 # everything else (paddles, ball) just set to 1
19        return I.astype(np.float).ravel()
20
21    def _createModel(self):
22        """Create a NN with one hidden layer, define loss function for
23        policy gradient learning and RMSProp optimizer.
24        - self.probs is used to make prediction of an action by the agent
25        - self.opt is used for policy update iteration
26        """
27        initializer = tf.contrib.layers.xavier_initializer()
28
29        with tf.variable_scope('policy'):
30            hidden = tf.layers.dense(self.states_ph, 200, activation=tf.nn.relu
, kernel_initializer = initializer)
31            logits = tf.layers.dense(hidden, 1, activation=None,
kernel_initializer = initializer)
32
33            self.probs = tf.sigmoid(logits, name="sigmoid")
34
35            cross_entropy = tf.nn.sigmoid_cross_entropy_with_logits(
36                labels=self.actions_ph, logits=logits, name="
cross_entropy")
37            loss = tf.reduce_sum(tf.multiply(self.rewards_ph, cross_entropy,
name="rewards"))
38
39            lr=1e-3

```

```

40         decay_rate=0.99
41         self.opt = tf.train.RMSPropOptimizer(lr, decay=decay_rate).minimize(
42             loss)
43     def predict(self, state):
44         """Return the probability of an action given a state, i.e. the policy
45         """
46         probs = self.sess.run(self.probs,
47                                feed_dict={self.states_ph: state.reshape((-1, state.size))})
48         return probs
49     def update_model(self, states, actions, rewards):
50         """Perform gradient ascent step and update model weights"""
51         self.sess.run(self.opt, feed_dict={self.states_ph: states,
52                                             self.actions_ph: actions,
53                                             self.rewards_ph: rewards})

```

Listing 5: Policy Gradient Brain definition

The `Environment` class implements the Policy Gradient algorithm shown in Figure 3 in which the agent observes a state, performs an action according to the current policy, and stores a transition in memory. This is repeated for `BATCH_SIZE` times, where in the current configuration `BATCH_SIZE` is set to 10. After that the weights are updated, then the memory is emptied before repeating another iteration.

```

1  class Environment():
2      def __init__(self, game):
3          self.game = game
4          self.env = gym.make(game)
5
6      def run(self, agent, max_episodes=1000):
7          """Policy gradient learning algorithm"""
8          tf.reset_default_graph()
9          observation = self.env.reset()
10         prev_s = None
11         epr = []
12         step = 0 # weight update step
13         episode_number = BATCH_SIZE*step
14         while True:
15             cur_s = agent.brain.preprocess(observation)
16             s = cur_s - prev_s if prev_s is not None else np.zeros((agent.
stateCnt))
17             prev_s = cur_s
18             # perform action and observe environment
19             action_label = agent.action_label(s)
20             action = action_label+2
21             observation, reward, done, info = self.env.step(action)
22             # append to memory
23             agent.memory['states'].append(s)
24             agent.memory['action_labels'].append(action_label)
25             epr.append(reward)
26             # end of one episode/game
27             if done:
28                 episode_number += 1
29                 discounted_epr = agent.discount_rewards(epr)
30                 agent.memory['rewards'] += discounted_epr
31                 R = sum(epr)
32                 epr = []
33                 if episode_number > max_episodes:
34                     break
35                 if episode_number % BATCH_SIZE == 0:

```

```

36         step += 1
37         # policy gradient weight update
38         agent.learn()
39         observation = self.env.reset()
40         prev_s = None

```

Listing 6: Policy Gradient Brain definition

3.3 Refinement

As previously mentioned, during initial testing the same CNN architecture used in [2] was adopted for both the algorithms, but afterwards this was abandoned as it revealed too computational expensive to test the codes in reasonable time. In addition, many implementations of the Policy Gradient algorithm that can be found for instance in GitHub use a discount reward calculated for the entire episode. Nevertheless, it was found that this approach prevented the agent to learn. The solution that enable to agent to learn was to adopt Karpathy’s interpretation of the discounted reward for Pong game [8], in which the discounted reward is calculated for a specific game and reset afterwards during an episode.

4 Results

The results of the agent performance during the episodes played are shown for the DQN and the Policy Gradient respectively in Figure 6a and 6b. These should be compared to the random agent benchmark shown in Figure 4. The DQN version of the agent was not successful in learning an optimal value-function during the 2000 episodes played, which took about 33 hours of computational time. Since it is reported that in DQN the performance can vary significantly in different training sessions, the algorithm was run twice but the results were similar. Increasing the exploration phase of the ϵ -greedy strategy was also not successful.

On the other hand, the Policy Gradient version started to learn after 500 episodes and was able to reach a return of approximately -14 after 2000 episodes. This required approximately 3.5 hours of computational time. The reason why the policy gradient is 10 times faster than the DQN is because it performs and update after a batch of episodes rather then after every time step. Therefore, the Policy Gradient has shown superior performance than the Policy Gradient and in a much shorter time

5 Conclusion

5.1 Reflection and improvement

Reinforcement learning is a very fascinating field of machine learning, that has recently attracted a lot of attention from the AI community due to the different problems where it could be applied with possibly superior performances of human beings. Although RL algorithms has been developed several years ago, significant results have been obtained only recently by combining them with neural network models. Despite the complexity of understanding and applying these methods, the main difficulty in using deep-RL was that it requires very long computational

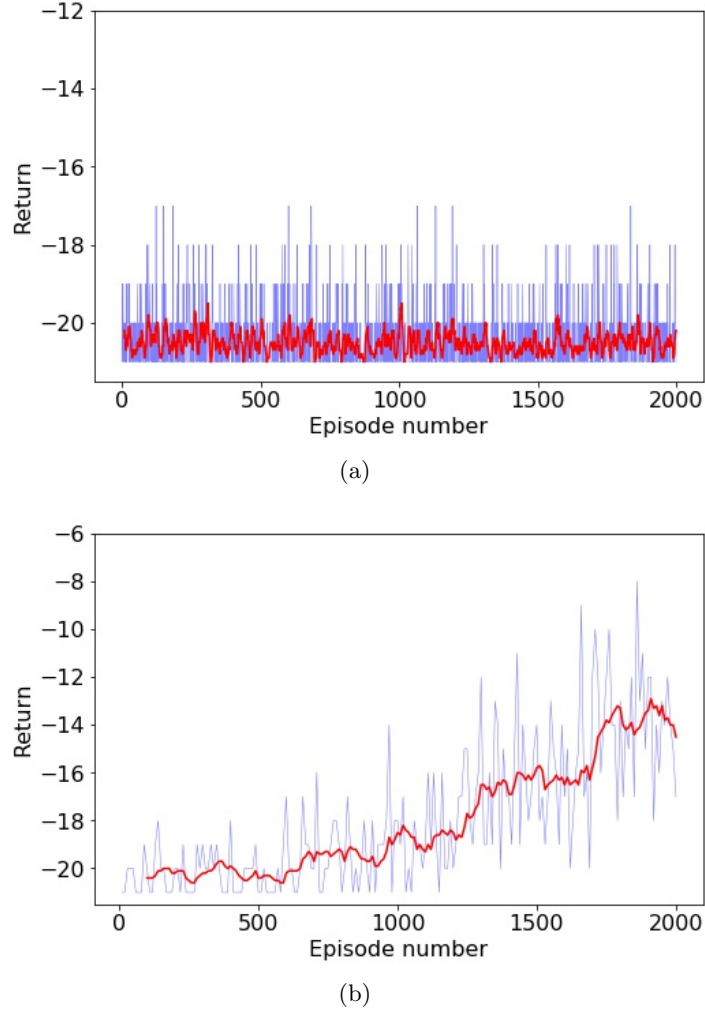


Figure 6: Agent training performance for (a) DQN and (b) Policy Gradient algorithm. In blue the return obtained at the end of an episode, and in red the moving average calculated in a window of 10 episodes.

times. This makes it difficult to test and debug the code when the algorithm is not performing well.

An obvious improvement to the current work would be to understand why the DQN algorithm is not performing better than a random agent after a relatively long training time. It would be interesting also to test more advanced versions of the algorithms such as Actor-Critic DQN, and a Policy Gradient version that uses the value function together with the discounted rewards [3].

References

- [1] V. Mnih et al. Playing Atari With Deep Reinforcement Learning. *NIPS Deep Learning Workshop*, 2013

- [2] V. Mnih et al. Human-level control through deep reinforcement learning. *Nature* 518 (7540):529–533, 2013
- [3] V. Mnih et al. Asynchronous Methods for Deep Reinforcement Learning. *International Conference on Machine Learning*, 2016
- [4] R.S. Sutton and A.G. Barto. *Reinforcement Learning: An Introduction*. Second Edition. Complete Draft, 2017
- [5] <http://gym.openai.com/envs/Pong-v0/>
- [6] <https://jaromiru.com/2016/10/03/lets-make-a-dqn-implementation/>
- [7] R.J. Williams. Simple statistical gradient-following algorithms for connectionist
- [8] <http://karpathy.github.io/2016/05/31/rl/>
- [9] <https://becominghuman.ai/lets-build-an-atari-ai-part-1-dqn-df57e8ff3b26>