

Reinforcement Learning

Kjell Raaijmakers (1244095), Jeroen van Riel (1236068)

June 26, 2022

Introduction

This report considers the problem of using reinforcement learning for training agents to play emulated games. We will first discuss a very simple toy game to explain and compare different learning algorithms. Next, we use the well-known game *Pong* to illustrate why using these reinforcement learning techniques may be hard on more complex games. Finally, we will present our efforts at implementing a reinforcement learning agent for the game *Breakout*.

Toy Game (Task 1)

We choose to explore the *Frozen Lake*¹ game, which involves a small rectangular grid world $W = \{F, H, S, G\}^{w \times h}$ containing w times h patches. Most patches are either frozen F or contain a hole H . There are two special patches S and G that indicate the *start* and *goal*, respectively, of an agent that moves in this grid world. Starting in S , the agent can choose to move a single step in one of the directions *up*, *right*, *down* or *left*, so we define the action space as $\mathcal{A} = \{U, R, D, L\}$. Due to the slippery nature of the frozen patches, the agent may not always end up in the chosen direction. The observation available to the agent is the position it actually ended up. When the agent ends up in a hole, the episode ends without any reward. If the agent is able to reach the goal patch G , the episode ends with a unit reward. Any other intermediate actions do not yield any reward.

Agents for Frozen Lake (Task 2 and 3)

Now that we have explained the game we will use the agents on, we will implement three learning agents on the game. First of all, we will look at a learning agent based on Monte Carlo (MC) prediction. For MC, our goal is to find each value for $q_\pi(s, a)$, where $q_\pi(s, a)$ is the expected return when starting in state s and taking action a , and then following policy π . The way the agent does this is by first generating an episode under a policy (in our case this will be either when the agent ends in a hole, or when he reaches the goal), and from there look for each first occurrence of being in state s and taking action a , and average the reward to the estimated value we already had of $q_\pi(s, a)$ (for example, if in 5 different episodes, we take action 0 in state 0 and get a reward of 1 in 1 of the episodes (and no rewards in the other episodes), our estimated value of $q_\pi(0, 0) = \frac{\text{Total rewards}}{\text{Total number of runs}} = \frac{1}{5}$). The agent also updates their strategy based on its current estimation of the matrix $Q = (q(s, a)_{s \in S, a \in A})$ (where S is the state space and A is the action space). In particular, the agents strategy of choosing an action will always be ε -greedy, i.e., the agent takes the optimal action based on their earlier observations in $1 - \varepsilon$ cases, and in the remaining ε cases, it takes a random action (so that the agent keeps exploring). The exact form of the algorithm is based on Figure 5.4 of [4]. In Figure 1 we see the success rate of the Monte-Carlo agent. Here, we decrease the ε over time. As we expect, we see that the success rate increases, which is normal since we decrease ε , thus decreasing the probability of taking a non-optimal action. However, as we will later see, this is not the best agent we can get.

¹https://www.gymnasium.ml/environments/toy_text/frozen_lake/

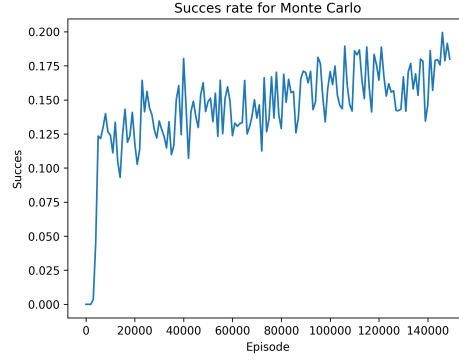
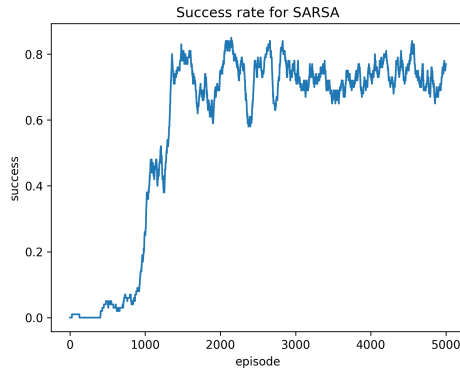


Figure 1: Success rate of the Monte-Carlo agent

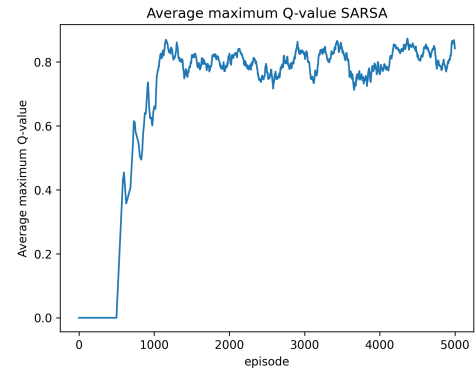
Then we will take a short look at temporal-difference (TD) learning. TD-learning, just like Monte Carlo, is a way to evaluate the values of being in a state and taking an action. However, for Monte Carlo, the agent waits for the whole run to be over, while TD-learning agents only look at the next step. This means that TD-learning agents may update matrix Q for each step, while Monte Carlo agents update the matrix for multiple steps only once an episode is finished.

Now, TD-learning only evaluates a given strategy π , however, we need something to update this strategy. We look at two different methods: SARSA and Q-learning. Looking at SARSA, after taking a step from state s with action a and observing reward R , we look in which state we would end up (which we denote by s_n) and which action a_n we would take under our current strategy. We use this information, as well as the reward we observe, to update our matrix of Q . The algorithm used for this algorithm can be found in Figure 6.9 of [4].

For SARSA, we see some of the results in Figure 2. First, in Figure 2a, we have the success rate for the SARSA agent, for decreasing values for ε . Interesting to see is that the rate of success is much higher than the results of the Monte-Carlo agent. Furthermore, we also looked at the maximum value of the Q matrix, averaged over the last 1000 values. We see that, after about 1000 episodes, this value begins to stabilize, which is logical, since we do not expect the Q -matrix to change as much when a lot of episodes are finished.



(a) Success rate for the SARSA agent



(b) Maximum Q -matrix values for the SARSA agent

Figure 2: Results for the SARSA agent

Finally, we look at Q-learning. Q-learning looks similar to SARSA, apart from the fact that we do

not take a new action a_n from the new state s_n , but instead look at the best action an_Q^* we can take under our current matrix Q . The algorithm used for this algorithm can be found in Figure 6.12 of [4].

Then for the results, in Figure 3, we see similar results as for the SARSA agent. This, can be explained by the fact that, when decreasing the value for ε , the Q -learning agent and the SARSA-agent look the same.

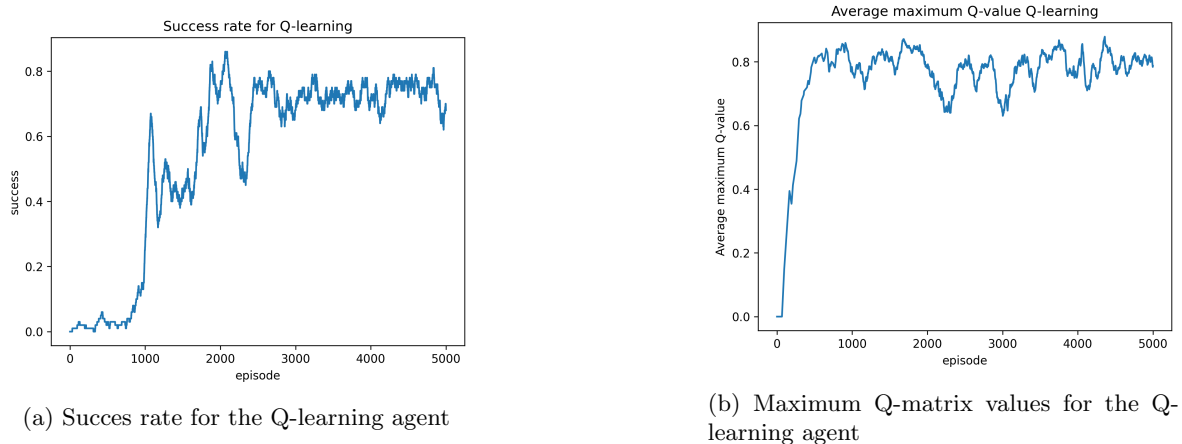


Figure 3: Results for the Q-learning agent

Before moving on to some more complex games, one might ask the question: why did the Monte-Carlo agent perform so much worse than the other agents? This might be due to the fact that we might not have done enough simulations, since we would expect that, when the number of episodes increases, the performance also increases.

Pong (Task 4)

Now that we have looked at a smaller game, we take a look at some more complex games. For the first game, which is the game Pong, we will observe the difficulties we can face when dealing with games which have a large amount of observations. After that (from Task 5 and onwards) we will report on the implementation of the learning agents we have looked at in Task 2 for the game of Breakout.

For the game of Pong, we are dealing with a image of size 210×160 , with color values which can range from 0 to 255.

This makes that the total observation space of Pong is big; an upperbound to the total number of observations is

$$210 \cdot 160 \cdot 256 = 8601600 \approx 8.6 \cdot 10^6$$

This is only considering the number of observations according to a single frame. However, this might not even be enough.

The agent would not be able to discriminate between situations with different velocity/acceleration from single frames alone. Since the agent has to map each state to an action at every time step, we need to incorporate information about velocity in the state representation. Assume for the moment that we have a method to extract the exact coordinates of the ball from each single frame. Taking a simple numerical approximation of velocity, as well as an approximation of the direction of the ball, we would be able to calculate where exactly we need to be at what time.

Now we can ask ourselves what the most important pieces of informations are, as taking $2.6 \cdot 10^7$ observations into account (without even looking at the speed of the ball), one might ask themselves if all this information is necessary. When looking at a single frame of the game (take for example

Figure 1 of the assignment), the most important pieces of information is the position of the ball and the position of the player. Furthermore, in the frame, a small trail can be seen, which can be used to find the direction of the ball. On top of that, we might not care that much about what happens close to the opponent; a small approximation error close to the opponent might lead to quite a large error along the way anyway. Then, since the color of the ball is constant, instead of checking all colors, we might only care about one color (or to reduce errors, a small array of colors). Finally, we still need something to deal with the speed of the ball, with which we deal after explaining our proposal of the reduced observation space.

Speaking of our reduced observation space, we might consider the following: Divide the opponents half in c_{opponent} even columns and r_{opponent} even rows, and divide the own half in c_{own} even columns and r_{own} even rows (where $c_{\text{opponent}} < c_{\text{own}}$ and $r_{\text{opponent}} < r_{\text{own}}$). Then, select an array of k colors, with the array of colors being close to the color of the ball (and by extend, the color of the own player, which is the same color). Finally, to deal with the velocity of the ball, we can take the last f frames, so that we can estimate the velocity of the ball.

An upperbound for this reduced state space is now

$$(c_{\text{opponent}} \cdot r_{\text{opponent}} + c_{\text{own}} \cdot r_{\text{own}}) \cdot k \cdot f$$

If we want to get a value for this, we can for example take $c_{\text{opponent}} = r_{\text{opponent}} = 4$, $c_{\text{own}} = r_{\text{own}} = 8$, $k = 10$ and $f = 5$, so that the upperbound is

$$(4 \cdot 4 + 8 \cdot 8) \cdot 10 \cdot 5 = 4000 = 4 \cdot 10^3$$

Given that our upperbound for the full state space was of the order 10^6 , this reduced state space is a lot smaller.

Playing Atari Games

The Arcade Learning Environment (ALE) [1] provides a collection of game emulators on which reinforcement learning algorithms can be tested. We will consider the *Breakout*² game in which the player has to control a paddle at the bottom of the screen in order to keep bouncing a ball towards the top of the screen where rows of bricks are stacked, see Figure 4. When the ball hits a brick, it breaks and the player gets a reward. When the player fails to bounce the ball back up, it disappears out of the game screen and the player loses a *life*. Once all lives are used, the game ends.

We refer to a completed game as an *episode* consisting of a number of discrete *steps* at which the player observes a screen *frame* and has to select an *action*. Each frame is encoded as a tensor of shape (210, 160, 3), where the last dimension holds the intensity of each RGB color channels. The four available actions are *do nothing* (0), *fire* (1), move the paddle *right* (2) and *left* (3). The fire action starts the game and uses a life to restart after failure. The number of steps per episode lies around 100 when the agent follows a completely random policy. Our goal is to train an agent that achieves a high total reward per episode.

There are multiple approaches we could take to tackle the stated problem. As illustrated in Section , the complexity of the observation space is a main issue in practical reinforcement learning. One could try to manually engineering usefull feature for the agent, such as velocity and acceleration of the ball. A major problem with this approach is that it is does not transfer easily to other problem domains, where the observations are of a completely different nature. Alternatively, deep learning methods have been proven useful in automatically constructing low-dimensional representation of the observations.

Deep Q-Learning

We will consider the use of a Deep Q-Network (DQN) [3], which is a Q-learning method that uses a Convolutional Neural Net (CNN) as a function approximator. Let $R_t = \sum_i^T \gamma^{i-t} r_i$ be the *discounted*

²<https://www.gymnasium.ml/environments/atari/breakout/>

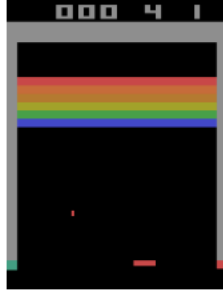


Figure 4: Original observation frame (without preprocessing) from a game of Breakout. The score indicated at the top left is still zero and there are four lives still available.

return at time t . The idea of Q-learning [5] is to estimate the optimal action-value function, which is defined as the expected discounted reward

$$Q^*(s, a) = \max_{\pi} \mathbb{E}[R_t | s_t = s, a_t = a, \pi]$$

that the agent receives when taking action a from state s and following policy π afterwards. We know that Q^* satisfies the *Bellman equation*

$$Q^*(s, a) = \mathbb{E}_{s'}[r + \gamma \max_{a'} Q^*(s', a') | s, a],$$

where the expectation is taken over next state s' , which is randomly drawn from the environment. It has been shown that iterating over

$$Q(s, a) \leftarrow \mathbb{E}_{s'}[r + \gamma \max_{a'} Q(s', a') | s, a]$$

converges to the optimal action-value function.

It is common to represent the action-value $Q(s, a; \theta) \approx Q^*$ using a function approximator to enable generalization between states. Without this method, the agent has to estimate $Q^*(s, a)$ for each distinct state-action pair, which is very restrictive in the current setting because of the large number of such pairs. Like in the original papers, we use a convolutional neural network (CNN) with three layers and parameters θ , which can be trained by optimizing in each iteration the loss function

$$L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho} [(y_i - Q(s, a; \theta_i))^2],$$

where $y_i = \mathbb{E}_{s'}[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$ is referred to as the *target* of the iteration. The expectation in the loss function is taken over the *behavioral distribution* ρ , which is due to the *behavioral policy* that is followed during training to collect samples (s, a, r, s') to estimate both expected values. We used a simple ϵ -greedy policy π_ϵ based on the current estimate of the state-action function. Because Q-learning estimates the action-values belonging to the optimal policy π^* while following policy π_ϵ , it is referred to as an *off-policy* method.

Like the original paper, we used a *experience replay* mechanism to smooth the learning. Samples (s, a, r, s') are stored in fixed-length memory, where old samples are removed when the memory is full and a new sample is added. During each iteration, the algorithm samples a random batch from memory and performs a gradient update. This ensures that the behavioral policy does not change too fast, which smoothes the learning and prevents oscillations.

Training

We used PyTorch to implement the model. In order to reduce training time, we applied several preprocessing steps. The Gym package for Python provides an easy way to *wrap* existing environments

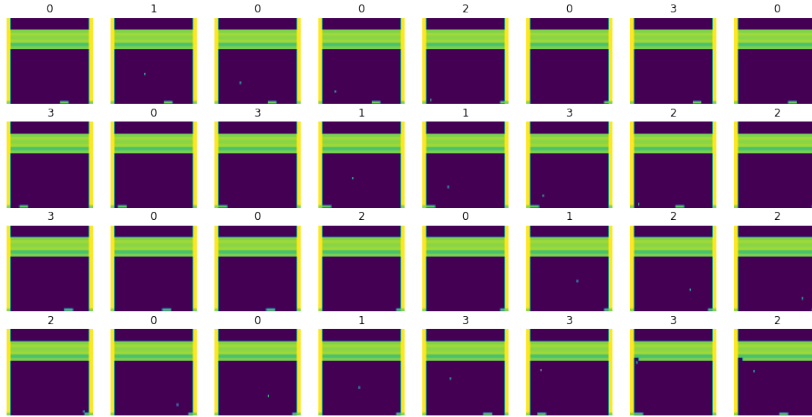


Figure 5: Some game preprocessed observations with $k = 4$ (so 3 frames are skipped). Note that the pixels have only one channel (grayscale), but we used a blue to yellow colorscheme for a clearer picture. The action number is indicated above the observation that resulted from it. As you can see, the game starts after the player has used the *fire* action.

to alter the interaction with an agent ³. We implemented a wrapper to apply (i) image preprocessing, (ii) frame skipping, (iii) frame stacking and (iv) automatic restarts.

Each screen frame includes a the score and the number of remaining lives, which is not essential for solving the problem at hand. After the top and bottom part of a frame is removed, we rescale it to 84×84 and merge the 3 color channels by taking the mean value, see Figure 5

Like the original paper, we added *frame skipping* to speed up the training process. Upon receiving an action from the agent, our wrapper environment submits the same action to the emulator for $k - 1$ steps, while accumulating rewards. The resulting state and accumulated reward is returned, such that the agent effectively observes every k th step. Another wrapper implements *frame stacking*, which involves providing the last l frames to the agent. This is a common technique to allow the agent to learn about velocity and acceleration. For a detailed explanation of how the authors of the original paper on DQN implemented both steps, we refer to [2].

One important issue that initially prevented the agent from learning anything within reasonable time is the fact that the fire action is required to start the game and to restart after the ball has been missed. Initially we considered the fire action as one of the actions the agent must choose from. Later we decided to automate the (re)starting, to speed up the learning process. The observations from the Gym environment include a number that indicates the number of lives left. Based on this number, the wrapper decides to do the fire action. The DQN agent now only has to learn how to control the paddle through the other three actions.

Evaluation

Initial development and testing of the implementation was done using Google Colab ⁴, but their free tier restricts use of GPU environments for extended time. Therefore, we activated free \$100 Microsoft Azure ⁵ credit that is in the Github Student Developer Pack⁶. Using this, we were able to train and evaluate our model in Azure Machine Learning Studio ⁷, which also provides the popular NVIDIA

³<https://www.gymnasium.ml/content/wrappers/>

⁴<https://colab.research.google.com>

⁵<https://azure.microsoft.com>

⁶<https://education.github.com/pack>

⁷<https://azure.microsoft.com/en-us/services/machine-learning/>

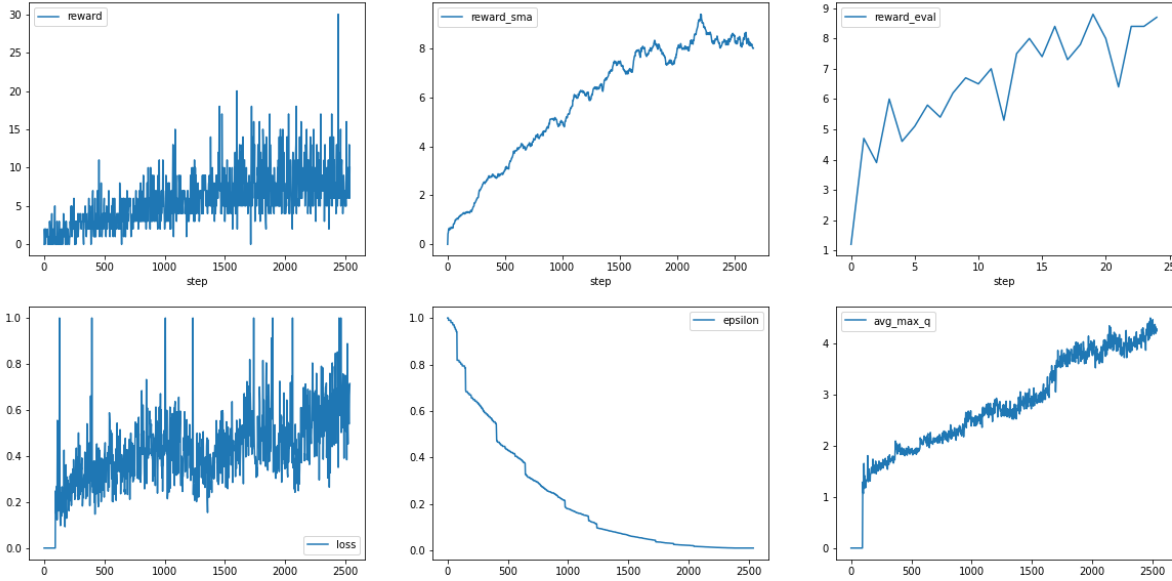


Figure 6: Evaluation of our final implementation.

Tesla K80 ⁸ that was available in Colab. Our source code and training logs can be found in our online repository ⁹. For easy reference, the source code has also been included in Appendix .

We performed various experiments, while changing parameters for the memory size, number of frame skips and learning rate. We did observe oscillating learning behavior with small memory size of 10000, but we encountered also a situation in which a memory size of only 1000 steps clearly showed a growing learning curve. Unfortunately, we had not enough time to develop a systematic approach to have a clear comparison of results, so we showcase one training run here

We trained the model for $T = 2500$ episodes, using $\gamma = 0.95$, learning rate $\alpha = 0.0005$ and memory size of $M = 100000$ steps. We started with totally random exploration $\epsilon = 1$ and decrease after each 1000 steps using $\epsilon = 0.99^{t/1000}$. The training progress is shown in Figure 6. The top row shows various measures of the reward: reward per episode, simple moving average and the third plot shows the periodic evaluation of a 0.05-greedy policy using the trained Q-values, averaged over 10 episodes. We do not plot discounted reward, because we are eventually interested in the total accumulated reward. The bottom row shows various other measures that we used to ensure smooth learning. It is important that the loss does not show too large oscillations, which may indicate a too large learning rate. Finally, the last plot in the bottom row shows the average maximum Q values in a minibatch, where the maximum is taken over the actions.

For the purpose of illustration and to help during debugging, we recorded a video of an 0.05-greedy agent after each 100 episodes. We had to go through a lot of Gym source code before we found a little bug that prevented us from using the provided `RecordVideo` wrapper. The bug was related to the way some metadata property was named, which prevented the recording wrapper to silently fail.

During training there appeared two episodes that are particularly long. We are sure that this can only be due to the fire action not being provided, but that would mean there is a bug in our wrapper, which we have not yet been able to find. Unfortunately, it does affect the learning, because ϵ is decreased with the number of steps. This explains the large drops in the plot for ϵ in Figure 6.

⁸<https://www.nvidia.com/en-gb/data-center/tesla-k80/>

⁹<https://github.com/jeroenvanriel/decision-theory/tree/master/assignment-3>

Ideas

To account for the sparsity of rewards, we planned to apply some curriculum engineering based on the principle of *optimism in face of uncertainty* by forcing the algorithm to collect more samples with positive reward. Instead of storing all the samples that we encounter in the collection phase, we randomly skip a fraction of samples that have zero reward.

We tried using double Q learning, but lack of computation time prevented us from producing a comparison.

References

- [1] M. G. Bellemare et al. “The Arcade Learning Environment: An Evaluation Platform for General Agents”. In: *Journal of Artificial Intelligence Research* 47 (June 2013), pp. 253–279. DOI: 10.1613/jair.3912. URL: <https://doi.org/10.1613/jair.3912>.
- [2] *Frame Skipping and Pre-Processing for Deep Q-Networks on Atari 2600 Games*. <https://danieltakeshi.github.io/2016/11/25/frame-skipping-and-preprocessing-for-deep-q-networks-on-atari-2600-games/>. Accessed: 2022-06-26.
- [3] Volodymyr Mnih et al. *Playing Atari with Deep Reinforcement Learning*. 2013. DOI: 10.48550/ARXIV.1312.5602. URL: <https://arxiv.org/abs/1312.5602>.
- [4] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. Second. The MIT Press, 2018. URL: <http://incompleteideas.net/book/the-book-2nd.html>.
- [5] Christopher J.C.H. Watkins and Peter Dayan. “Technical Note: Q-Learning”. In: *Machine Learning* 8.3 (May 1992), pp. 279–292. ISSN: 1573-0565. DOI: 10.1023/A:1022676722315. URL: <https://doi.org/10.1023/A:1022676722315>.

Appendices

Tasks

Jeroen van Riel typed Task 1, 5, 6 and 7 (onward from “Playing Atari games”) and made the code for these chapters as well, as well as helping with the code for Task 2,3,4. Kjell Raaijmakers typed Task 2, 3, 4, and made the code for Task 2,3,4.

Code Monte-Carlo Frozen Lake

```
1 import gym
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import random as rand
5 from queue import Queue
6 from random import random, choice
7
8 env = gym.make('FrozenLake-v1', desc=["SFFF", "FHFH", "FFFH", "HFFG"])
9
10
11 class EpsilonSoftPolicy:
12     def __init__(self, epsilon, env):
13         self.epsilon = epsilon
14         self.obs_space = range(env.observation_space.n)
15         self.action_space = range(env.action_space.n)
16         self.actions = {
17             state : env.action_space.sample()
18             for state in self.obs_space
19         }
20
21     def act(self, state):
```



```

22         if random() < self.epsilon:
23             # choose random action
24             return choice(self.action_space)
25         else:
26             # choose greedy action
27             return self.actions[state]
28
29     def improve(self, state, a_star):
30         self.actions[state] = a_star
31
32     def print(self):
33         print(self.actions)
34
35     def change_epsilon(self, new_epsilon):
36         self.epsilon = new_epsilon
37
38
39     def generate_episode(policy):
40         state_actions, rewards = [], []
41         t = 0
42         state = env.reset()
43         while True:
44             action = policy.act(state)
45             state_actions.append((state, action))
46
47             t += 1
48             state, reward, done, info = env.step(action)
49             rewards.append(reward)
50
51             if done:
52                 break
53
54         return state_actions, rewards
55
56     class SMA:
57         """Simple moving average with incremental update."""
58         def __init__(self, k):
59             self.q = Queue(k)
60             self.k = k
61             self.SMA = 0
62
63         def put(self, p):
64             if self.q.full():
65                 self.SMA += (p - self.q.get()) / self.k
66             else:
67                 self.SMA += p / self.k
68             self.q.put(p)
69
70         def get(self):
71             return self.SMA
72
73
74
75
76     # on-policy Monte Carlo control
77

```

```

78
79 gamma = 0.95
80 sma = SMA(1000)
81
82 epsilon_all = [0.1]
83 nriters = 150
84 finalrewards = np.zeros((len(epsilon_all),nriters))
85
86
87 for x in range(len(epsilon_all)):
88     epsilon = epsilon_all[x]
89     policy = EpsilonSoftPolicy(epsilon, env)
90     n = 0
91     Q = np.empty((env.observation_space.n, env.action_space.n))
92     N = np.ones((env.observation_space.n, env.action_space.n))
93     T=0
94     G = 0
95     while T<nriters:
96         policy.change_epsilon(max(2*epsilon-(1/75)*epsilon*T,0))
97         sa, rewards = generate_episode(policy)
98         for t in range(len(sa) - 1, -1, -1):
99             G = gamma * G + rewards[t]
100             if sa[t] not in sa[0:t]:
101                 # incremental update of state action value
102                 Q[sa[t]] = Q[sa[t]] + (G - Q[sa[t]]) / N[sa[t]]
103                 N[sa[t]] += 1
104
105                 # improvement step
106                 state, action = sa[t]
107                 a_star = np.argmax(Q[state,:])
108                 policy.improve(state, a_star)
109
110             sma.put(G)
111             if n >= 1000:
112                 n = 0
113                 print(x,T)
114                 finalrewards[x][T] = sma.get()
115                 T=T+1
116             else:
117                 n += 1
118
119
120 x_plot=range(nriters)
121 x_plot = np.multiply(x_plot, 1000)
122
123 plt.plot(x_plot,finalrewards[0])
124 plt.title("Success rate for Monte Carlo")
125 plt.xlabel("Episode")
126 plt.ylabel("Success")
127 plt.savefig('Monte_frozen_lake_changing_epsilon_1.png', dpi=300, bbox_inches='
tight')
128 plt.show

```

Code SARSA Frozen Lake

```

1 import gym
2 import numpy as np

```

```

3 import matplotlib.pyplot as plt
4 import math
5 from queue import Queue
6 from random import random, choice
7
8
9 class EpsilonSoftPolicy:
10     def __init__(self, epsilon, env):
11         self.epsilon = epsilon
12         self.obs_space = range(env.observation_space.n)
13         self.action_space = range(env.action_space.n)
14         self.actions = {
15             state : env.action_space.sample()
16             for state in self.obs_space
17         }
18
19     def act(self, state):
20         if random() < self.epsilon:
21             # choose random action
22             return choice(self.action_space)
23         else:
24             # choose greedy action
25             return self.actions[state]
26
27     def improve(self, state, a_star):
28         self.actions[state] = a_star
29
30     def print(self):
31         print(self.actions)
32
33     def change_epsilon(self, new_epsilon):
34         self.epsilon = new_epsilon
35
36 class SMA:
37     """Simple moving average with incremental update."""
38     def __init__(self, k):
39         self.q = Queue(k)
40         self.k = k
41         self.SMA = 0
42
43     def put(self, p):
44         if self.q.full():
45             self.SMA += (p - self.q.get()) / self.k
46         else:
47             self.SMA += p / self.k
48         self.q.put(p)
49
50     def get(self):
51         return self.SMA
52
53
54 env = gym.make('FrozenLake-v1', desc=["SFFF", "FHFH", "FFFH", "HFFG"])
55
56 Q = np.zeros((env.observation_space.n, env.action_space.n))
57
58

```

```

59 gamma = 0.95
60 alpha = 0.5
61
62 sma = SMA(1000)
63
64 epsilon_all = [0.2]
65 niterations = 250
66 finalrewards = np.zeros((len(epsilon_all),niterations))
67
68 for x in range(len(epsilon_all)):
69     epsilon = epsilon_all[x]
70     policy = EpsilonSoftPolicy(epsilon, env)
71     n = 0
72     T=0
73     G = 0
74     while T<niterations:
75         policy.change_epsilon(max(epsilon-(1/niterations)*epsilon*T,0))
76         s = env.reset()
77         while True:
78             a = policy.act(s)
79             sn, reward, done, info = env.step(a)
80             G = gamma * G + reward
81             an = policy.act(sn)
82
83             Q[s,a] = Q[s,a] + alpha * (reward + gamma * Q[sn,an] - Q[s,a])
84
85             # update greedy action
86             a_star = np.argmax(Q[s,:])
87             policy.improve(s, a_star)
88
89             s = sn
90             a = an
91             if done:
92                 sma.put(G)
93                 if n >= 1000:
94                     n = 0
95                     finalrewards[x][T] = sma.get()
96                     print(T)
97                     T=T+1
98                 else:
99                     n += 1
100             break
101
102
103
104 x_plot=range(niterations)
105 plt.plot(x_plot,finalrewards[0])
106 plt.title("Success rate for SARSA")
107 plt.xlabel("Episode")
108 plt.ylabel("Success")
109 plt.savefig('SARSA_frozen_lake_changing_epsilon.png', dpi=300, bbox_inches='
    tight')
110 plt.show

```

Q-learning Frozen Lake

```

1 import gym

```

```

2 import numpy as np
3 import matplotlib.pyplot as plt
4 from queue import Queue
5 from random import random, choice
6
7
8 class EpsilonSoftPolicy:
9     def __init__(self, epsilon, env):
10         self.epsilon = epsilon
11         self.obs_space = range(env.observation_space.n)
12         self.action_space = range(env.action_space.n)
13         self.actions = {
14             state : env.action_space.sample()
15             for state in self.obs_space
16         }
17
18     def act(self, state):
19         if random() < self.epsilon:
20             # choose random action
21             return choice(self.action_space)
22         else:
23             # choose greedy action
24             return self.actions[state]
25
26     def improve(self, state, a_star):
27         self.actions[state] = a_star
28
29     def print(self):
30         print(self.actions)
31
32 class SMA:
33     """Simple moving average with incremental update."""
34     def __init__(self, k):
35         self.q = Queue(k)
36         self.k = k
37         self.SMA = 0
38
39     def put(self, p):
40         if self.q.full():
41             self.SMA += (p - self.q.get()) / self.k
42         else:
43             self.SMA += p / self.k
44         self.q.put(p)
45
46     def get(self):
47         return self.SMA
48
49
50 env = gym.make('FrozenLake-v1', desc=["SFFF", "FHFH", "FFFH", "HFFG"])
51
52 Q = np.zeros((env.observation_space.n, env.action_space.n))
53 gamma = 0.9
54 alpha = 0.5
55 episodes = 1000
56 T = 10000 # maximum number of steps in one episode
57

```

```

58 sma = SMA(100)
59 episode_sma = np.zeros((episodes))
60
61 def update_epsilon(epsilon):
62     #return min(epsilon, max(0, 2 * epsilon - (1 / 50) * epsilon * episodes))
63     return epsilon - 0.001
64
65 policy = EpsilonSoftPolicy(1, env)
66 for episode in range(episodes):
67     total_reward = 0
68     s = env.reset()
69     for t in range(T):
70         a = policy.act(s)
71         sn, reward, done, info = env.step(a)
72         total_reward += reward
73
74         Q[s,a] = Q[s,a] + alpha * (reward + gamma * np.max(Q[sn,:]) - Q[s,a])
75
76         # update greedy action
77         a_star = np.argmax(Q[s,:])
78         policy.improve(s, a_star)
79
80         s = sn
81
82         if done:
83             sma.put(total_reward)
84             episode_sma[episode] = sma.get()
85             policy.epsilon = update_epsilon(policy.epsilon)
86             print(f'epsilon: {policy.epsilon}, episode: {episode}, sma: {sma.
get()}, steps: {t}')
87             break
88
89 x_plot = range(episodes)
90 plt.plot(x_plot, episode_sma)
91 plt.title('Success rate for Q-learning')
92 plt.xlabel('episode')
93 plt.ylabel('success')
94 plt.savefig('q-frozen-lake.png', dpi=300, bbox_inches='tight')
95 plt.show

```

Atari Breakout

train.py

```

1 import gym, torch, cv2, random
2 import numpy as np
3 from collections import deque
4 from datetime import datetime
5 from gym.wrappers import TransformObservation, FrameStack, RecordVideo
6
7 from gym import logger
8 logger.set_level(logger.INFO)
9
10 # for logging in azure machine learning studio
11 from azureml.core import Run
12 run = Run.get_context()
13

```

```

14 DEV = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
15 print(f'using device {DEV}')
16
17 def get_time():
18     now = datetime.now()
19     return now.strftime("%H:%M:%S")
20
21 def preprocess(frame):
22     frame = frame[34:-16, :, :] # crop
23     frame = cv2.resize(frame, (84, 84))
24     frame = frame.mean(-1) # to grayscale
25     frame = frame.astype('float') / 255 # scale pixels
26     return frame
27
28 class FrameSkip(gym.Wrapper):
29     def __init__(self, env, frames):
30         super().__init__(env)
31         self.env = env
32         self.frames = frames
33         self.lives = 5
34
35     def step(self, action):
36         action = {0:0, 1:2, 2:3}[action] # map to NOPE, LEFT, RIGHT
37         reward = 0
38         for _ in range(self.frames):
39             s, r, done, info = self.env.step(action)
40             lives_next = info['lives']
41             reward += r
42             if done:
43                 break
44
45             if self.lives != lives_next and not done:
46                 s, r, done, info = self.env.step(1) # FIRE to restart
47
48         return s, reward, done, info
49
50     def reset(self):
51         super().reset()
52         s, _, _, _ = self.env.step(1) # FIRE to start
53         return s
54
55 def wrap(env):
56     env = TransformObservation(env, preprocess)
57     env = FrameSkip(env, 2)
58     env = FrameStack(env, num_stack=4)
59     return env
60
61 env = wrap(gym.make('ALE/Breakout-v5'))
62
63 # maximum number of emulation steps
64 T_MAX = 100000
65
66 # evaluation is done in a different environment, with recording enabled
67 EVAL_INTERVAL = 100
68 EVAL_EPISODES = 10
69 env_eval = gym.make('ALE/Breakout-v5')

```

```

70 env_eval.metadata['render.modes'] = env_eval.metadata.get('render.modes', [])
    # fix a gym bug
71 vid = RecordVideo(env_eval, 'outputs/', episode_trigger=lambda: True)
72 env_eval = wrap(vid)
73
74 def evaluate(agent, episode):
75     epsilon = agent.epsilon
76     agent.epsilon = 0.05
77
78     rewards = []
79     for t in range(EVAL_EPISODES):
80         if t == 0: # big hack to only record first with the episode number
81             vid.episode_id = episode
82             vid.episode_trigger = lambda x: True
83         else:
84             vid.episode_trigger = lambda x: False
85         s = env_eval.reset()
86         reward = 0
87         for _ in range(T_MAX):
88             a = agent.get_action(s)
89             s, r, done, _ = env_eval.step(a)
90             reward += r
91             if done:
92                 break
93         rewards.append(reward)
94
95     agent.epsilon = epsilon # reset back
96     return np.mean(rewards)
97
98 if __name__ == "__main__":
99     from agent import Agent
100     EPISODES = 10000
101     MEMORY_SIZE = 100000 # steps
102     MIN_MEMORY = 10000 # steps
103     bs = 64
104     lr = 5e-5
105     gamma = 0.95
106     agent = Agent(env, lr, gamma, MEMORY_SIZE, MIN_MEMORY, bs, DEV)
107
108     reward_sma = deque(maxlen=100) # simple moving average
109
110     total_step = 0 # number of steps over all episodes
111     for episode in range(EPISODES):
112         s = env.reset()
113
114         total_reward = 0
115         total_loss = 0
116         total_max_q = 0
117         for t in range(T_MAX): # max number of actual emulation steps
118             a = agent.get_action(s)
119             s_next, r, done, info = env.step(a)
120             agent.store(s, a, r, s_next, done)
121             s = s_next
122
123             loss, max_q = agent.train()
124

```



```

125         total_loss += loss
126         total_max_q += max_q
127         total_reward += r
128
129         total_step += 1
130         if total_step % 1000 == 0:
131             agent.epsilon = max(0.01, agent.epsilon * 0.99)
132
133         if done: # end of episode
134             agent.update_target_dqn()
135
136             # periodic evaluation and model saving
137             if total_step > MIN_MEMORY and episode % EVAL_INTERVAL == 0:
138                 torch.save(agent.online_dqn, f'outputs/online-{episode}.pt
139
140                 # torch.save(agent.target_dqn, f'outputs/target-{episode}.
141                 pt')
142
143                 print('model saved!')
144
145                 reward_eval = evaluate(agent, episode)
146                 print(f'reward_eval: {reward_eval}')
147                 run.log('reward_eval', reward_eval)
148
149             # logging
150             reward_sma.append(total_reward)
151             avg_max_q = total_max_q / t
152             run.log('reward', total_reward)
153             run.log('reward_sma', np.mean(reward_sma))
154             run.log('loss', total_loss)
155             run.log('avg_max_q', avg_max_q)
156             run.log('epsilon', agent.epsilon)
157             print(f'[{get_time()}] episode: {episode} reward: {int(
total_reward)} reward_sma: {np.mean(reward_sma):.3f} loss: {total_loss:.3f
} avg_max_q: {avg_max_q:.3f} epsilon: {agent.epsilon:.2f} last_step: {t}
total_step: {total_step}')
158             break

```

agent.py

```

1  import torch
2  import random
3  import numpy as np
4  from torch.optim import Adam
5  from torch.nn.functional import mse_loss, smooth_l1_loss
6  from collections import deque
7
8  from model import DQN, DQN1, DQN3
9
10 def q_for_action(qs, actions):
11     """Get the q values from (64, 3) corresponding to the action."""
12     return qs.gather(1, actions.unsqueeze(1)).squeeze(1)
13
14 class Agent:
15     def __init__(self, env, lr, gamma, memory_size, min_memory, batch_size,
16         device):
17         self.epsilon = 1

```

```

18     self.dev = device
19     self.online_dqn = DQN(3).to(device).train()
20     self.target_dqn = DQN(3).to(device)
21     self.target_dqn.load_state_dict(self.online_dqn.state_dict())
22     self.target_dqn.eval()
23
24     self.gamma = gamma
25
26     self.memory = deque(maxlen=memory_size)
27     self.batch_size = batch_size
28     self.min_memory = min_memory
29
30     self.optimizer = Adam(self.online_dqn.parameters(), lr=lr)
31
32     self.opt_exp_count = 0
33     self.just_exp_count = 0
34
35     def store(self, s, a, r, s_next, done):
36         self.memory.append([s, a, r, s_next, done])
37
38         # optimistic experience
39         # if len(self.memory) > self.min_memory:
40         #     self.memory.append([s, a, r, s_next, done])
41         #     print('just store')
42         #     self.just_exp_count += 1
43         # else:
44         #     if r > 0:
45         #         self.memory.append([s, a, r, s_next, done])
46         #         self.opt_exp_count += 1
47         #     elif random.random() < 0.5:
48         #         self.memory.append([s, a, r, s_next, done])
49         #         self.just_exp_count += 1
50
51     def get_action(self, state):
52         if random.random() < self.epsilon:
53             # choose random action
54             return random.choice(range(3))
55         else:
56             # return greedy action
57             x = torch.unsqueeze(torch.from_numpy(np.asarray(state)), dim=0).
float().to(self.dev)
58             return torch.argmax(self.online_dqn(x)).item()
59
60     def double_q_loss(self, s_b, a_b, r_b, s_next_b, done_b):
61         # get q values for current states
62         q_vals = self.online_dqn(s_b)
63         # actual q values for selected actions
64         actual_q_vals = q_for_action(q_vals, a_b)
65
66         # target q values
67         # TODO: maybe no_grad() is required here!
68         next_q_vals = self.online_dqn(s_next_b) # shape: N x 4
69         next_actions = next_q_vals.max(1)[1]
70         next_target_q_vals = q_for_action(self.target_dqn(s_next_b),
next_actions)
71

```

```

72     # double Q learning target
73     target = r_b + self.gamma * next_target_q_vals * (1 - done_b)
74
75     return mse_loss(actual_q_vals, target.detach()), q_vals
76
77     def td_loss(self, s_b, a_b, r_b, s_next_b, done_b):
78         q_vals = self.online_dqn(s_b)
79         # select q values corresponding to actions
80         q_vals = torch.squeeze(torch.take_along_dim(q_vals, a_b.unsqueeze(1),
1))
81
82         with torch.no_grad():
83             qnext = self.online_dqn(s_next_b)
84             m, _ = torch.max(qnext, dim=-1)
85             target = r_b + self.gamma * m * (1 - done_b)
86
87         # return mse_loss(q_vals, target), q_vals
88         return smooth_l1_loss(q_vals, target), q_vals
89
90     def train(self):
91         if len(self.memory) < self.min_memory:
92             return 0, 0
93
94         # sample minibatch
95         s_b, a_b, r_b, s_next_b, done_b = zip(*random.sample(self.memory, self
.batch_size))
96
97         s_b = torch.from_numpy(np.asarray(s_b)).float().to(self.dev)
98         a_b = torch.tensor(a_b, dtype=torch.long, device=self.dev)
99         r_b = torch.tensor(r_b, dtype=torch.float, device=self.dev)
100        s_next_b = torch.from_numpy(np.asarray(s_next_b)).float().to(self.dev)
101        done_b = torch.tensor(done_b, dtype=torch.float, device=self.dev)
102
103        self.optimizer.zero_grad()
104        loss, q_vals = self.td_loss(s_b, a_b, r_b, s_next_b, done_b)
105        loss.backward()
106        self.optimizer.step()
107
108        return loss.item(), q_vals.max().item()
109
110    def update_target_dqn(self):
111        self.target_dqn.load_state_dict(self.online_dqn.state_dict())
112        self.target_dqn.eval()

```

model.py

```

1  import torch.nn as nn
2
3  class DQN(nn.Module):
4      def __init__(self, n_actions):
5          super().__init__()
6          self.n_actions = n_actions
7          self.model = nn.Sequential(
8              # 4 input frames
9              nn.Conv2d(4, 32, 8, stride=4),
10             nn.ReLU(),
11             nn.Conv2d(32, 64, 4, stride=2),

```

```

12         nn.ReLU(),
13         nn.Conv2d(64, 64, 3, stride=1),
14         nn.ReLU(),
15         nn.Conv2d(64, 1024, 7, stride=1),
16         nn.ReLU(),
17         nn.Flatten(),
18         nn.Linear(1024, n_actions),
19     )
20
21     def forward(self, x):
22         return self.model(x)
23
24
25     class DQN1(nn.Module):
26         def __init__(self, n_actions):
27             super().__init__()
28             self.n_actions = n_actions
29             self.model = nn.Sequential(
30                 # 4 input frames
31                 nn.Conv2d(4, 32, 8, stride=4),
32                 nn.ReLU(),
33                 nn.Conv2d(32, 64, 4, stride=2),
34                 nn.ReLU(),
35                 nn.Conv2d(64, 64, 3, stride=1),
36                 nn.ReLU(),
37                 nn.Flatten(),
38                 nn.Linear(7 * 7 * 64, 512),
39                 nn.ReLU(),
40                 nn.Linear(512, n_actions),
41             )
42             # for i in [0, 2, 4, 7, 9]:
43             #     nn.init.xavier_uniform(self.model[i].weight)
44
45         def forward(self, x):
46             return self.model(x)
47
48     class DQN2(nn.Module):
49         def __init__(self, n_actions):
50             super().__init__()
51             self.n_actions = n_actions
52             self.model = nn.Sequential(
53                 # 4 input frames
54                 nn.Conv2d(4, 32, 8, stride=4),
55                 nn.BatchNorm2d(32),
56                 nn.ReLU(),
57                 nn.Conv2d(32, 64, 4, stride=2),
58                 nn.BatchNorm2d(64),
59                 nn.ReLU(),
60                 nn.Conv2d(64, 64, 3, stride=1),
61                 nn.BatchNorm2d(64),
62                 nn.ReLU(),
63                 nn.Flatten(),
64                 nn.LazyLinear(256),
65                 nn.ReLU(),
66                 nn.Linear(256, n_actions),
67             )

```

```

68
69     def forward(self, x):
70         return self.model(x)
71
72 class DQN3(nn.Module):
73     def __init__(self, n_actions):
74         super().__init__()
75         self.n_actions = n_actions
76         self.model = nn.Sequential(
77             # 4 input frames (84, 84)
78             nn.Conv2d(4, 32, 8),
79             nn.MaxPool2d(2, stride=2),
80             nn.ReLU(),
81             nn.Conv2d(32, 64, 4),
82             nn.MaxPool2d(2, stride=2),
83             nn.ReLU(),
84             nn.Conv2d(64, 64, 3),
85             nn.MaxPool2d(2, stride=2),
86             nn.ReLU(),
87             nn.Flatten(),
88             nn.Linear(7 * 7 * 64, 512),
89             nn.ReLU(),
90             nn.Linear(512, n_actions),
91         )
92
93     def forward(self, x):
94         return self.model(x)

```