

# Autonomously Learning Systems WS21/22

## Assignment 2 Deep Reinforcement Learning

Ozan Özdenizci, Horst Petschenig

Presentation:	17.12.2021 10:00
Info Hour:	07.01.2022 10:00, Cisco WebEx meeting, see TC
Deadline:	<b>17.01.2022 08:00 (extended deadline)</b>
Hand-in procedure:	Use the <b>cover sheet</b> from the TeachCenter Submit your <b>Python files and report</b> on the TeachCenter
Course info:	<a href="https://tc.tugraz.at/main/course/view.php?id=3110">https://tc.tugraz.at/main/course/view.php?id=3110</a>
Group size:	up to two students

### General remarks

Your submission will be graded based on:

- Correctness (Is your code doing what it should be doing?)
- The depth of your interpretations (Usually, only a couple of lines are needed)
- The quality of your plots (Is everything clearly visible in the print-out? Are axes labeled? ...)
- Your submission should run with Python 3.7+
- Both the submitted report and the code. If some results or plots are not described in the report, they will **not** count towards your points!
- Code in comments will not be executed, inspected or graded. Make sure that your code is runnable.

## 1 Deep Q-Learning for Atari Breakout [5 Points]

We will solve the Breakout environment from the OpenAI Gym Framework, using off-policy Q-learning with non-linear function approximation via deep neural networks. We will particularly use the deep Q-learning (DQN) framework, as introduced by [Mnih et al., 2015], which exploits two important mechanisms, namely a *fixed target Q* and *experience replay*. We will use PyTorch for this example, such that we can harness automatically computed gradients for the parameters to be optimized.

Open the file `deep_q_network.ipynb` code skeleton which you can run on Google Colab with GPUs. This file creates the BreakoutNoFrameskip-v4 environment. Breakout is the well-known Atari game environment where the player moves a board at the bottom of the screen that returns a ball to destroy the blocks at the top of the screen (see Figure 1). To solve the Q function learning problem, while one can use a very large set inputs to the agent that represent the environment state, we will use deep convolutional neural networks which can solve the same problem by only looking at the 2D scene frames.

- a) **Using a Fixed Target Q Network:** Parameter updates for standard Q-learning (with function approximation) after taking an action  $a_t$  in state  $s_t$  and resulting in the state  $s_{t+1}$  by observing the reward  $r_{t+1}$  is  $\theta_{t+1} \leftarrow \theta_t + \alpha[y_t^Q - Q_\theta(s_t, a_t)]\nabla_\theta Q_\theta(s_t, a_t)$ , where  $y_t^Q = r_{t+1} + \gamma \max_{a'} Q_\theta(s_{t+1}, a')$ . Essentially this is interpreted as stochastic gradient descent optimization by updating the  $Q_\theta(s_t, a_t)$  towards  $y_t^Q$  using a mean-squared error loss. To overcome the instability in optimization that can be caused when small updates to  $Q_\theta(s_t, a_t)$  change the policy and potentially the target data distribution, DQNs implement a fixed target Q network as an important component.

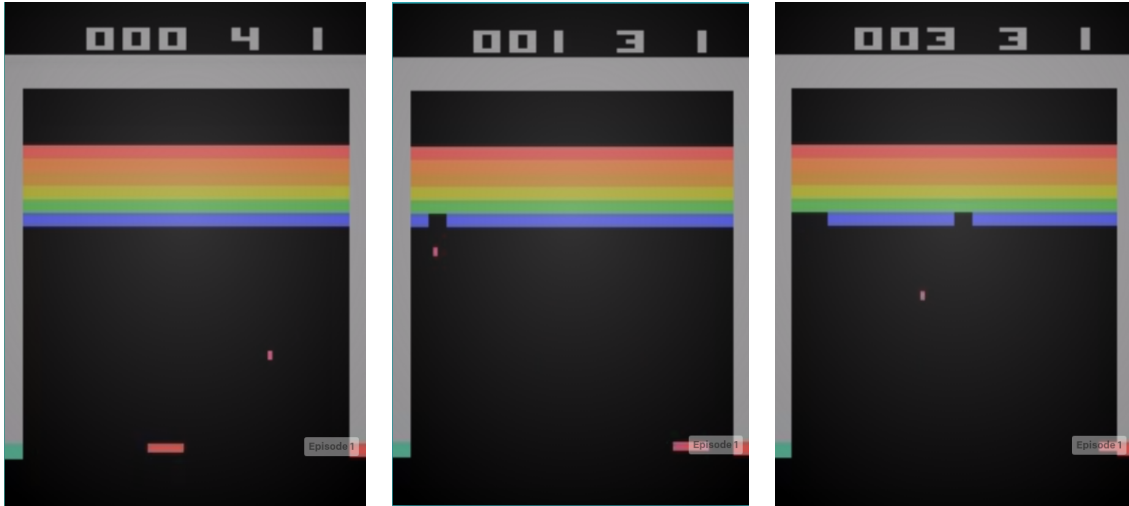


Figure 1: An illustration of the Breakout Atari environment.

The fixed target Q network with parameters  $\bar{\theta}$  has the identical convolutional architecture as the online Q network, except that we will not update the parameters  $\bar{\theta}$  at each step but copy the parameters from the online Q network after every  $k$  steps (by setting  $\bar{\theta} = \theta$ ), and keep it fixed otherwise. In this case  $y_t^{\text{DQN}}$  for stochastic gradient descent optimization will be written as:

$$y_t^{\text{DQN}} = r_{t+1} + \gamma \max_{a'} Q_{\bar{\theta}}(s_{t+1}, a'), \quad (1)$$

resulting in the following parameter updates with mean-squared error loss:

$$\theta_{t+1} \leftarrow \theta_t + \alpha [r_{t+1} + \gamma \max_a Q_{\bar{\theta}}(s_{t+1}, a) - Q_{\theta}(s_t, a_t)] \nabla_{\theta} Q_{\theta}(s_t, a_t). \quad (2)$$

Fill the TODOs in the provided code skeleton to implement a convolutional network for the DQN (i.e., explore different configurations by varying the number of Conv2D layers, activations or batch normalization), and solving the task with a fixed target and an online DQN with periodic updates to the target model. Explain and explore what is the functionality of the variable `burn_in_phase`. You will also see some utility functions to be completed (or modified if you wish) that will pre-process input frames to reduce the input complexity by converting them to grayscale and downsampling of original frames with dimensionality  $210 \times 160$  into frames of  $84 \times 84$ .

- b) **Experience Replay Memory:** We will now implement an experience replay memory by storing state transitions  $(s_t, a_t, r_{t+1}, s_{t+1})$  that the agent observes, which allows us to reuse this data by uniformly sampling from the buffer  $\mathcal{B}$  later, i.e.,  $(s_t, a_t, r_{t+1}, s_{t+1}) \sim U(\mathcal{B})$ . This framework overcomes the problem of repeatedly observing correlated data in our sequence of observations. We will implement a simple replay memory with a double ended queue, to sample an uncorrelated minibatch of transitions from this cyclic buffer before the parameter optimization steps.

Fill the TODOs in the experience replay memory function to perform appropriate buffering and uniform sampling while training your models consistently through the rest of the pipeline.

- c) **Optimization Configurations:** We conventionally used a mean-squared error loss. Explore the use of a Huber loss<sup>1</sup> function. Explain mathematically what is the main difference between these two? Then implement these two separately and investigate how different do they perform. Using the model you implemented, observe the influence of initialization for the parameter  $\epsilon$ . Report the accumulated reward obtained in the test episodes with initial values for  $\epsilon \in \{0, 0.1, 0.2, 0.5, 1\}$ . What is the best value? Is it still possible to solve the task with  $\epsilon = 0$  and why?

<sup>1</sup><https://pytorch.org/docs/stable/generated/torch.nn.HuberLoss.html>

## 2 Combining Improvements for DQNs [5+2\* Points]

We will now sequentially incorporate improvements to the baseline DQN model we just implemented. For each model and additional improvement mechanism that will be implemented, make sure to report your results with plots where you present the increasing obtained accumulated reward by the agent (on the y-axis) across increasing number of frames (on the x-axis) in the test episodes.

- a) **Double Q-Learning:** This approach proposed by [van Hasselt et al., 2016] exploits two Q functions to disentangle the max operation in Eq. (2), where one model will be used to determine the greedy action selection at  $s_{t+1}$  and another one to evaluate this selected action. This will eventually help reducing the overestimation of Q values and allow better and faster learning. To combine this idea with DQNs, we can use our target Q network with parameters  $\bar{\theta}$  (which we have been periodically updating the weights of) as a proxy for the second Q function.

We will decompose the max operation in Eq. (1) where the action is both selected and evaluated by the target Q network, and instead perform action selection using the online Q network with parameters  $\theta$ , and evaluate this action on the target Q network. Hence,  $y_t^{\text{DDQN}}$  will be:

$$y_t^{\text{DDQN}} = r_{t+1} + \gamma Q_{\bar{\theta}}(s_{t+1}, \arg \max_{a'} Q_{\theta}(s_{t+1}, a')). \quad (3)$$

Create a boolean variable `run_as_ddqn` at the top of your script which would make your model run as a Double DQN when `True`. Compare performance of DDQN with respect to the DQN model.

*Hint:* You will only need to modify the loss function computation with the new  $y_t^{\text{DDQN}}$ . You can simply implement this by performing the argmax action selection for  $s_{t+1}$  with the online Q network, and then gather that selected action's output value from the target Q network for  $s_{t+1}$ .

- b) **Multi-step Learning:** Previously we have been accumulating a single reward  $r_{t+1}$  and then choosing a greedy action at the next step. Alternatively, we can use multiple steps to accumulate rewards. The truncated  $n$ -step return from a given state  $s_t$  is given by

$$r_t^{(n)} = \sum_{k=0}^{n-1} \gamma^k r_{t+k+1}. \quad (4)$$

Note that for  $n = 1$  this corresponds to the standard  $r + 1$  single reward setting. Accordingly, implement a multi-step learning variant for your DQN by adapting  $y_t^{\text{DQN}}$  into:

$$y_t^{\text{n-DQN}} = r_t^{(n)} + \gamma^{(n)} \max_{a'} Q_{\bar{\theta}}(s_{t+n}, a'), \quad (5)$$

which can be incorporated into our usual mean-squared error or Huber loss function to optimize  $Q_{\theta}(s_t, a_t)$ , as well as also combined with a DDQN via Eq. (3). Perform multi-step learning with  $n = 3$  steps. Compare the  $n$ -step learning approach with respect to a baseline DQN. Incorporate this addition to the DDQN as well (i.e.,  $n$ -step DDQN) and report comparisons.

*Hint:* For the  $n$ -step learning variant you may want to adapt the experience replay function too.

- c) **[Bonus: 2 Points] Prioritized Experience Replay:** Your implementation of DQN so far uniformly samples experiences from the replay buffer. It would make sense however to sample more frequently those transitions from which there is much to learn. Implement a separate prioritized experience sampling function for your DQN, which samples transitions proportional to a probability relative to the last encountered absolute temporal difference error, i.e.,  $|r_{t+1} + \gamma \max_{a'} Q_{\bar{\theta}}(s_{t+1}, a') - Q_{\theta}(s_t, a_t)|$ . Perform prioritized experience replay memory also for the DDQN models where the absolute temporal difference is then defined as  $|r_{t+1} + \gamma Q_{\bar{\theta}}(s_{t+1}, \arg \max_{a'} Q_{\theta}(s_{t+1}, a')) - Q_{\theta}(s_t, a_t)|$ . Report performance with prioritized replay in comparison to the baseline DQN and DDQN.