

## Assignment 5: Q-Learning, Double Q-learning, and DQN

**Due: Apr. 21th, 2024, 23:59pm**

**Question 1.** Consider a reinforcement-learning system with two states (namely,  $s_1$  and  $s_2$ ) and two actions (namely,  $a_1$  and  $a_2$ ). Suppose that a Q-learning trial has been conducted with the agent transitioning through the following state sequence by taking the actions as indicated below:

$$s_1 \xrightarrow{a_1, 1} s_1 \xrightarrow{a_2, 1} s_2 \xrightarrow{a_1, 10} s_1 \xrightarrow{a_2, 1} s_2 \quad (1)$$

where the number following the action above an arrow indicates the reward received by the agent upon taking that action. For instance, the first arrow implies that the agent at state  $s_1$  takes action  $a_1$ , which results in the agent remaining in state  $s_1$  and receiving a reward of 1. Complete the table below to show the values of the Q-function at the end of each action taken by the agent during the trial. For instance, the value of  $Q(s_1, a_1)$  is to be entered in the top left empty cell in the table shown. Assume that the initial values of the Q-learning are 0. Use a fixed learning rate of  $\alpha = 0.5$  and a discount rate of  $\gamma = 0.5$ .

Q	$s_1$	$s_2$
$a_1$		
$a_2$		

**Note:** Show your detailed calculation steps for obtaining these Q-function values. There are four actions for this trial, so your answer should include four such tables, one for each action taken.

**Question 2 (Programming).** Consider the grid world shown in the following figure. This is a standard undiscounted, episodic task, with start and goal states, and the usual actions causing movement up, down, right, and left. If the agent takes an action that would move it off the grid, it remains in its current position instead. The reward is -1 on all transitions except those into the black region. Stepping into the black region incurs a reward of -100 and sends the agent instantly back to the start.

S					G

- Implement Q-learning and SARSA on this task respectively with probability for exploration  $\epsilon = 0.1$ , step size  $\alpha = 0.1$ , and discount factor  $\lambda = 1$ . Choose the number of episodes sufficiently large (e.g., 500) so that a stable policy is learned.
  - Plot a figure with two curves that show the “Sum of rewards during the episode” against “Episodes” for Q-learning and SARSA respectively.
  - Plot the learned policy for each method.

Note that a script is given for your convenience of programming, so you can focus more on coding the Q-function update (search for **TODO** in the code source file to locate where you need to complete the code). To have a smoother plot, the given script performs 100 independent runs.

- Consider now running both Q-learning and SARSA for 500 training episodes first using  $\alpha = 0.1$ ,  $\gamma = 1$ , and  $\epsilon$ -greedy policy for collecting samples where  $\epsilon = 0.1$ . Then, take the greedy policy resulting from the learned Q-functions of both methods and run another 200 episodes following the obtained policy (that is, stop training after 500 episodes, but test the learned policies from both methods afterward).

- Plot the cumulative rewards for each episode for both SARSA and Q-Learning.

Note that

- Since the greedy policy obtained from the learned Q-function is deterministic, there is no need to perform independent runs to smooth the rewards curves anymore after 500 episodes.
  - It is possible that if you train the agent for not enough number of episodes (say the training episodes = 100), the learned policy could not even lead the agent to the Goal.
- Change  $\epsilon = \frac{1}{10t+1}$ , where  $t$  is the number of episodes for SARSA. Choose the total number episodes of training to be 500. Work on Task 2 again. Show the policy learned by SARSA and Q-learning. Plot the cumulative rewards for each episode for both SARSA and Q-learning.

In your report, include all the required plots for all 3 tasks. Comment on your simulation results. Submit your codes together with the report.

**Question 3 (Programming).** Implement Deep Q-learning Network (DQN) and DDQN (Double DQN) to balance an inverted **pendulum** from the OpenAI Gym environment. In this task, there is only one major incomplete code file **dqn.py** (Some miscellaneous local settings and customization of the code might be needed). Upon completion, one should be able to run the **train.py** file and successfully control the inverted pendulum to keep the upright position after training. All other codes are complete.

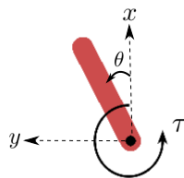


Figure 1: Balancing a pendulum using DQN, DDQN

### Tasks:

- Fill in the blanks in the **dqn.py** file marked with **TODO** to develop the following agents
  - DQN** agent, where Q-network takes 1 image and the angular velocity (i.e., the state  $s$ ), and the torque (i.e., the action  $a$ ) as input, and outputs the value  $Q(s, a) \in \mathbb{R}$  at the current  $(s, a)$  pair.
  - DQN** agent, where Q-network takes 4 consecutive images (i.e., the state  $s$ ) and the torque (i.e., the action  $a$ ) as input, and outputs the value  $Q(s, a) \in \mathbb{R}$  at the current  $(s, a)$  pair.
  - DDQN** agent, where Q-network takes 4 consecutive images as input (i.e., state  $s$ ), and outputs the values at all actions, that is, outputs a vector  $Q(s, \cdot) \in \mathbb{R}^{|A|}$  and each element of the vector corresponding to  $Q(s, a)$  at an action  $a$ .
- Run the **train.py** file to train a model such that it can keep the inverted pendulum upright. Record a video to show that the inverted pendulum is upright under the developed 3 different agents. Submit the completed code and the video.

3. Plot a figure that shows the “mean\_ep\_100\_return” against the “used step” as shown in `train.py` for each agent. Play with the parameters in `train.py` file to see how it will affect the return curves. Comment on your exploration/discovery in your accompanying report.
4. In the `train.py` file, we use an 9-dimensional vector to denote the discrete action space. Explore more high-dimensional vectors to see whether they can lead to better learning results.

### Technical Tips (More in the Appendix 0.1):

- You can refer to this [Git repo](#) for reference, based on which we adapted our code. The link is also included in the comments of the given code. To complete the `learn()` function, you can refer to their implementation at [here](#). (For CNN implementation in Pytorch, there are a lot of examples online, you can also refer to their implementation [here](#) as well). You can refer to other code samples, as it is fairly standard to train DNN in PyTorch.
- The observation from the environment we provided in the code consists of the latest four consecutive frames of images and angular velocities, i.e., `observation = [image, velocity]`, where `dim(image) = 4x42x42` and `dim(velocity) = 1 × 4`.
- Stacking of 4 most recent images and angular velocities is set in `wrappers.py`. Specifically, in the `class FrameStack`.
- The higher the index is, the newer the data is. E.g., from index 3 to 0, the corresponding data is from the most recent to the oldest.
- Note that in preprocessing, the RGB channels of the image are combined into one channel in the code. See `class ProcessFrame42` in `wrappers.py`:

```
1 img = img[:, :, 0] * 0.299 + img[:, :, 1] * 0.587 + img[:, :, 2] * 0.114
```

- In the setup of the given code, the Q-network takes the observation and the action torque as the input and outputs a Q-value for this state-action pair (please refer to the `forward` function in `dqn.py`). You will need to modify the code to accommodate the different setups described in the problem requirement.

### Notes:

- Two example videos are given to you to show the initial output and the desired output for Question 3.
- One should submit a report, all completed codes, and an output video.
- **You can modify all the codes freely as long as you can complete all the tasks successfully.**
- Refer to the Appendix below for tips on coding environment setup, dependence installation, etc.

## Submission

- Make sure to submit one PDF report, one MP4 video for each type of the 3 NN agents, and all the PYTHON codes. Include all your files into a folder and compress it to a zip file. Only submit one zip file for this assignment.
- Naming convention: Assignment5\_YourMatricNumber\_YourName.zip
- **Submission deadline:** Apr. 21th, 2024, 23:59pm.

## Appendix

### Working Environment:

- We recommend using [PyCharm](#) or [VScode](#) with [Anaconda](#) as your development environment, either in Windows, Mac, or Linux.

### Installing Dependency:

- You can use anaconda to create a python3 environment by installing the required packages listed in `environment.yml`:

```
1 cd DQN_DIRECTORY
2 conda env create -f environment.yml
```

- If some error messages from Anaconda are raised, you could choose to install the required python3 package manually. Run the following command with CMD in Windows or Shell in Linux or MacOS:

```
1 pip3 install pygame gym==0.26.1 opencv_python
```

- Install the **CPU-only version** of PyTorch as follows (if you don't use GPU):

```
1 pip3 install torch
```

- Note that if you want to use **PyTorch with CUDA support** for GPU acceleration, you need to install the appropriate version of PyTorch that is compiled with CUDA support. You can visit [PyTorch Get Started](#) to choose your options and get the installation guidance. For example, for Windows or Linux, run the following command, as prompted in the website (see the snapshot below):

```
1 //note that in windows, you need to add the --user flag as done below
2 pip3 install torch torchvision torchaudio --index-url https://download.pytorch.org/whl/cu118
  --user
```

PyTorch Build	Stable (2.2.2)		Preview (Nightly)	
Your OS	Linux	Mac	Windows	
Package	Conda	Pip	LibTorch	Source
Language	Python		C++ / Java	
Compute Platform	CUDA 11.8	CUDA 12.1	ROCM 5.7	CPU
Run this Command:	<pre>pip3 install torch torchvision torchaudio --index-url https://download.pytorch.org/whl/cu118</pre>			

### Test Your Built Environment (Dependency):

- When testing the built environment, you could let the code idle by running the following command in the terminal:

```
1 cd DQN_DIRECTORY
2 python3 train.py --idling
```

- If there is no error, that means you have installed all dependencies successfully. You can proceed to fill in the blanks in the `dqn.py` file marked with `TODO`.

### How to use:

- After completing all blanks in `dqn.py`, you can run `train.py` either in Vscode/PyCharm or in terminal by the following commands

```
1 cd DQN_DIRECTORY
2 python3 train.py
```

- For more details, please refer to `readme.md` file in DQN folder.

## 0.1 Tips:

- The state contains 4 consecutive images of  $4 \times 42 \times 42$ , and 4 consecutive angular velocities, returned in the lines below. You can check it using line 3.

```
1 s = env.reset() # reset environment
2 s_, r, done, infos, _ = env.step(a) # take a step
3 print(s[0][0].shape) # images
4 print(s[0][1].shape) # angular velocities
```

- The code only stops training and begins to animate the pendulum and save the animation if the average return over the last 100 training episodes is above -200 (You can change it to -3000 to see and save the animation without any training). See the code below in `train.py`,

```
1 # render
2 if mean_100_ep_return >= -200:
3     evaluate_performance(dqn, disc_actions, episode_length=EPISODE_LENGTH)
```

where the `evaluate_performance()` function will call the OpenAI Gym `env.render()` API to animate the experiment. `evaluate_performance()` is set to run 200 steps by default.

- To enable multiple workers\parallel threadings to obtain samples from multiple environments, you need to generate the same number of actions according to the number of environments (the variable `N_ENVS` in `train.py`). Specifically, the function `choose_action()` should return an array of the length of `N_ENVS`. For example, the following line samples `N_ENVS` number of actions from `[0, N_ACTIONS)` randomly:

```
1 action = np.random.randint(0, N_ACTIONS, N_ENVS)
```

- The data is saved to binary files (with `.pkl` name extension) in the default folders `data/plots`. The data is only saved once every `SAVE_FREQ` steps in `train.py`:

```
1 # save frequency
2 SAVE_FREQ = int(1.5e+4) // N_ENVS
```

- I included the code for saving the animation in `train.py` as the `save_animation` function, which calls the OpenAI Gym video recorder API (which requires installing the package `moviepy`):

```
1 pip install moviepy
```

It depends on the following function,

```
1 from gym.wrappers.monitoring.video_recorder import VideoRecorder
```

- The convergence and performance of the training can be affected by the exploration scheme (that is, related to the epsilon of  $\epsilon$ -greedy policy, and how  $\epsilon$  diminishes), as well as your design convolutional neural networks. One hidden linear layer (in the fully connected layer of CNN) is typically enough. More layers will incur more parameters and slow down the training. If designed properly, it should converge within 1 hour of training.
- Returning only angular velocity is set in `Pendulum.py`:

```
1 def _get_obs(self):
2     theta, thetadot = self.state
3     return np.array([thetadot], dtype=np.float32) # return only angular velocity
```