

HW10: Reinforcement Learning

Due Apr 29 by 11:59pm **Points** 100 **Submitting** a file upload **File Types** zip
Available Apr 20 at 12am - Apr 29 at 11:59pm 10 days

This assignment was locked Apr 29 at 11:59pm.

Getting started

The starter code is available here: [starter.zip](#) ↓

(https://canvas.wisc.edu/courses/230450/files/19704119/download?download_frd=1)

You can create and activate your virtual environment with the following commands:

```
python3 -m venv /path/to/new/virtual/environment
```

```
source /path/to/new/virtual/environment/bin/activate
```

once you have sourced your environment, you can run the following commands to install the necessary dependencies:

```
pip install --upgrade pip
```

```
pip install torch==1.8.0+cpu torchvision==0.9.0+cpu torchaudio==0.8.0 -f https://download.pytorch.org/whl/torch\_stable.html \_\(https://download.pytorch.org/whl/torch\_stable.html\)
```

```
pip install gym==0.17.1 numpy==1.19.5
```

You should now have a virtual environment which is fully compatible with the skeleton code. You should set up your environment on an instructional machine to do your final testing. More detail about python virtual environments is available in the assignment 6 prompt: [HW6: Neural Networks](#) .

Q-Learning

For the Q-learning and SARSA portion of HW10, we will be using the environment FrozenLake-v0 from OpenAI gym. This is a discrete environment where the agent can move in the cardinal directions, but is not guaranteed to move in the direction it chooses. The agent gets a reward of 1 when it reaches the tile marked G, and a reward of 0 in all other settings. You can read more about FrozenLake-v0 here:

<https://gym.openai.com/envs/FrozenLake-v0/> [\(https://gym.openai.com/envs/FrozenLake-v0/\)](https://gym.openai.com/envs/FrozenLake-v0/)

You will not need to change any code outside of the area marked TODO, but you are free to change the hyper-parameters if you want to. The update rule for Q-learning is:

$$Q(s, a) = \begin{cases} Q(s, a) + \alpha(r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a)) & \text{if !done} \\ Q(s, a) + \alpha(r - Q(s, a)) & \text{if done} \end{cases}$$

$$Q(s, a) = \begin{cases} Q(s, a) + \alpha(r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a)) & \text{if !done} \\ Q(s, a) + \alpha(r - Q(s, a)) & \text{if done} \end{cases} \quad \text{for each sampled } (s, a, r, s', \text{done}) \text{ tuple}$$

The agent should act according to an epsilon-greedy policy as defined in the Reinforcement Learning 1 slides.

In this equation, alpha is the learning rate hyper-parameter, and gamma is the discount factor hyper-parameter.

HINT: tests.py is worth looking at to gain an understanding of how to use the OpenAI gym env.

For this section, you should submit the files Q_learning.py and Q_TABLE.pkl

SARSA

SARSA is very similar to Q-learning, although it differs in the fact that it is on-policy. The name stands for "State, Action, Reward, State, Action," and refers to the tuples that are used by the update rule. You will not need to change any code outside of the area marked TODO, although, you are again welcome to change any hyperparameters you want.

The update rules for SARSA is:

$$Q(s, a) = Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a)) \quad \text{always} \quad \text{for each sampled } (s, a, r, s', a', \text{done}) \text{ tuple}$$

$$Q(s', a') = Q(s', a') + \alpha(r' - Q(s', a')) \quad \text{if done}$$

Similarly to Q-learning, the agent should act according to an epsilon-greedy policy as defined in the Reinforcement Learning 1 slides.

HINT: tests.py is worth looking at to gain an understanding of how to use the OpenAI gym env.

HINT: make sure that the SARSA update is performed for the last action-reward pair for each episode, because for FrozenLake-v0 this is only possible time for a non-zero reward to be received. For the last State action pair, the expected cumulative reward is equal to the reward for that state and action, since there will be no further rewards.

For this section, you should submit the files SARSA.py and SARSA_Q_TABLE.pkl

DQN (Extra Credit)

Deep Q-learning has produced some exciting results in the past 5 years. First introduced in 2015 in [Mnih et al.](https://storage.googleapis.com/deepmind-media/dqn/DQNNaturePaper.pdf) [_ \(https://storage.googleapis.com/deepmind-media/dqn/DQNNaturePaper.pdf\)](https://storage.googleapis.com/deepmind-media/dqn/DQNNaturePaper.pdf), DQN allows the application of Q-learning to domains with continuous observation spaces.

For the experience replay buffer, you should add tuples of the form $(s, a, r, s', \text{done})$. The sample_minibatch function should return a list containing a random sample of MINI_BATCH_SIZE tuples from the experience replay buffer

For the first TODO in the main while loop, you should use the current policy_net policy to determine the best action to take then add the resulting $(s, a, r, s', \text{done})$ tuple you encounter from the environment to the experience replay buffer.

HINT: remember to use "with torch.no_grad():" when you do not need to calculate gradients.

HINT: [Mnih et al.](https://storage.googleapis.com/deepmind-media/dqn/DQNNaturePaper.pdf) [_ \(https://storage.googleapis.com/deepmind-media/dqn/DQNNaturePaper.pdf\)](https://storage.googleapis.com/deepmind-media/dqn/DQNNaturePaper.pdf) has the full pseudo-code, and is useful to look at for reference when implementing DQN.

For the second TODO, you will need to sample MINI_BATCH_SIZE tuples from the experience replay buffer, then calculate the appropriate estimates using the policy_net, and y values using the observed rewards, second states, and the target_policy_net.

OpenAI gym Environment

You will need to use several OpenAI gym functions in order to operate your gym environment for reinforcement learning. As stated in a previous hint, tests.py has a lot of the function calls you need. Several important functions are as follows:

env.step(action)

Given that the environment is in state s , step takes an integer specifying the chosen action, and returns a tuple of the form $(s', r, \text{done}, \text{info})$. 'Done' specifies whether or not s' is the final state for that

particular episode, and 'info' is unused in this assignment.

env.reset()

Resets the environment to it's initial state, and returns that state.

env.action_space.sample()

Samples and integer corresponding to a random choice of action in the environment's action space.

env.action_space.n

In the setting of the environments we will be working with for these assignments, this is an integer corresponding to the number of possible actions in the environment's action space.

You can read more about OpenAI gym here: <https://gym.openai.com/docs/>
(<https://gym.openai.com/docs/>).

Submission Format

You can test your learned policies, by calling `python3 ./tests.py`. Make sure to test your saved Q-tables using `tests.py` on the instructional machines with a virtual environment set up as specified above. This is the same program, which we will be using to test your Q-tables and Q-network, so you will have a good idea about how many points you will receive for the automated tests portion of the grade. Your submission should contain the following files:

Required: `Q_learning.py`, `Q_TABLE.pkl`, `SARSA.py`, `SARSA_Q_TABLE.pkl`

Optional: `DQN.py`, `DQN.mdl`, `DQN_DATA.pkl`

Please submit these files in a zipped folder title `<yournetid>.zip` , where 'yournetid' is your net ID. Please make sure that there is not a folder inside the zipped folder, and that the submitted files are at the top level of the zipped folder.

The assignment is due April 29th at 11:59pm central time. We are not accepting late submissions for this assignment. Regrades will only be accepted if they are due to an error in `tests.py`.

HW10: Reinforcement Learning

Processing math: 100%

Criteria	Ratings			Pts
Q-Learning: Trained Q table 100 episode average score	30 pts Full Marks 100 episode average >= 0.6	15 pts Partial marks 100 episode average >= 0.4	0 pts No Marks 100 episode average <0.4	30 pts
Q-Learning: Implementation correctness	20 pts Full Marks Correct implementation of Q-learning	10 pts Partial marks Partially correct implementation of Q- learning	0 pts No Marks Incorrect implementation of Q- learning	20 pts
SARSA: Trained Q table 100 episode average score	30 pts Full Marks 100 episode average >= 0.6	15 pts Partial marks 100 episode average >= 0.4	0 pts No Marks 100 episode average < 0.4	30 pts
SARSA: Implementation correctness	20 pts Full Marks Correct implementation of SARSA	10 pts Partial marks Partially correct implementation of SARSA	0 pts No Marks Incorrect implementation of SARSA	20 pts
DQN: Trained Q network 100 episode average score	0 pts 0 points extra credit	0 pts 10 points extra credit 100 episode average >= 70		0 pts
Total Points: 100				

Processing math: 100%