Due Dec 14, 2021 by 10:59am **Available** Dec 1, 2021 at 11:59pm - Dec 14, 2021 at 11:15am Points 100 **Submitting** a file upload File Types zip

This assignment was locked Dec 14, 2021 at 11:15am.

Getting started

pip install --upgrade pip

The starter code is available here: <u>starter.zip</u> ↓

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 torch==1.8.0+cpu torchvision==0.9.0+cpu torchaudio==0.8.0 -f 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 that is fully compatible with the skeleton code. You should set up this virtual environment on an instructional machine (CSL machines) 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 (up, down, left, right), but is not guaranteed to move in the direction it chooses (more in the FrozenLake-v0 doc). 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/ \Rightarrow

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 (but keep the random seed unchanged). The update rule for Q-learning is:

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

The agent should act according to an epsilon-greedy policy as defined in the Reinforcement Learning 2 slides. Epsilon-greedy policy: link

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

Pseudo Code

else: pull current-best action HINT: tests.py is worth looking at to gain an understanding of how to use the OpenAI gym env.

p = random()

pull random action

if $p < \epsilon$:

A sample run of python Q_learning.py

100 EPISODE AVERAGE REWARD: 0.01

For this section, you should submit the files Q_learning.py and Q_TABLE.pkl. The Q_TABLE.pkl will be generated automatically when you run your Q_learning.py.

PSILON: 0.9038873549665959 LAST 100 EPISODE AVERAGE REWARD: 0.0 EPSILON: 0.8178301806491574

LAST 100 EPISODE AVERAGE REWARD: 0.03 EPSILON: 0.7399663251239436 [[ALAST 100 EPISODE AVERAGE REWARD: 0.02 EPSILON: 0.6695157201007336 LAST 100 EPISODE AVERAGE REWARD: 0.03 EPSILON: 0.6057725659163237 LAST 100 EPISODE AVERAGE REWARD: 0.04 EPSILON: 0.548098260578011 LAST 100 EPISODE AVERAGE REWARD: 0.03 EPSILON: 0.4959150020176678 LAST 100 EPISODE AVERAGE REWARD: 0.02 EPSILON: 0.44869999946146477 AST 100 EPISODE AVERAGE REWARD: 0.04 EPSILON: 0.4059802359226587 LAST 100 EPISODE AVERAGE REWARD: 0.12 EPSILON: 0.36732772934619257 LAST 100 EPISODE AVERAGE REWARD: 0.12 EPSILON: 0.33235524492954527 AST 100 EPISODE AVERAGE REWARD: 0.12 EPSILON: 0.3007124156643058 [[ALAST 100 EPISODE AVERAGE REWARD: 0.09 EPSILON: 0.2720822322326576 LAST 100 EPISODE AVERAGE REWARD: 0.16 EPSILON: 0.2461778670932771 The first few print outputs are shown above. Two things are worth pointing out:

accumulated reward obtained in this episode into the deque.

2. EPSILON is decreasing, which means you need to update (decay) the hyperparameter variable EPSILON in your code (e.g. using EPSILON_DECAY). Note: matching the sample output above is not guaranteed for the correctness of your program, please use the tests.py to test your code.

1. The value of EPISODE AVERAGE REWARD is obtained from the deque variable episode_reward_record. So you need to update this variable in your code by appending the

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

Similarly to Q-learning, the agent should act according to an epsilon-greedy policy as defined in the Reinforcement Learning 2 slides.

LAST 100 EPISODE AVERAGE REWARD: 0.01

LAST 100 EPISODE AVERAGE REWARD: 0.03

LAST 100 EPISODE AVERAGE REWARD: 0.01

LAST 100 EPISODE AVERAGE REWARD: 0.03

LAST 100 EPISODE AVERAGE REWARD: 0.08

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

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

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

A sample run of python SARSA.py

HINT: make sure that the SARSA update is performed for the last action-reward pair for each episode, because for FrozenLake-v0 this is the only possible time for a non-zero

EPSILON: 0.9038873549665959 LAST 100 EPISODE AVERAGE REWARD: 0.0 EPSILON: 0.8178301806491574

LAST 100 EPISODE AVERAGE REWARD: 0.03 EPSILON: 0.6057725659163237 LAST 100 EPISODE AVERAGE REWARD: 0.09 EPSILON: 0.548098260578011 LAST 100 EPISODE AVERAGE REWARD: 0.02 EPSILON: 0.4959150020176678

EPSILON: 0.44869999946146477

EPSILON: 0.7399663251239436

EPSILON: 0.6695157201007336

EPSILON: 0.4059802359226587 LAST 100 EPISODE AVERAGE REWARD: 0.11 EPSILON: 0.36732772934619257 LAST 100 EPISODE AVERAGE REWARD: 0.1 DQN (Extra Credit) Deep Q-learning has produced some exciting results in the past 5 years. First introduced in 2015 in Mnih et al. 🔁, DQN allows the application of Q-learning to domains with continuous observation spaces.

MINI_BATCH_SIZE tuples from the experience replay buffer

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.

For the first TODO in the main while loop, you should use the current policy to determine the best action to take then add the resulting (s, a, r, s', done) tuple

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

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.

OpenAl 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

not s' is the final state for that particular episode, and 'info' is unused in this assignment.

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

HINT: Mnih et al. ⇒ has the full pseudo-code, and is useful to look at for reference when implementing DQN.

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, done, info). 'Done' specifies whether or

Resets the environment to its initial state, and returns that state. env.action_space.sample()

action space.

env.reset()

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

You can read more about OpenAI gym here: https://gym.openai.com/docs/ □.

an error in tests.py.

Submission Format

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

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:

the submitted files are at the top level of the zipped folder. The assignment is due December 14th at 10:59 am central time. We are not accepting late submissions for this assignment. Regrades will only be accepted if they are due to

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

Grading

- 60 points (out of 100) will be given if you passed the tests.py (by running python tests.py after you generate the pkls). You can ensure that by testing on your own.

- 40 points will be determined via manual inspection of the code. We will inspect whether you implement the learning algorithms correctly.

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

- 10 points (extra credits) will be rewarded if your DQN passes the tests.py.

Criteria

Q-Learning: Trained Q table 100 episode average score

HW10: Reinforcement Learning

15 pts

10 pts

Partial marks

100 episode average >= 0.4

Ratings

0 pts

0 pts

No Marks

100 episode average < 0.4

Pts

30 pts

Total Points: 100

Q-Learning: Implementation correctness 20 pts

30 pts

Full Marks

100 episode average >= 0.6

	Full Marks Correct implementation of Q- learning	Partial marks Partially correct implementation of Q- learning		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		O 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