

Georgia Institute of Technology
Faculty of Interactive Computing

CS8003– ASSIGNMENT #1
Deep Reinforcement Learning

Instructor: Prof. Animesh Garg

Due date: 2024/10/22

Name: _____

Student Number: _____

This writing assignment contains 14 pages (including this cover page) and 28 questions. Total of points is 131.

Good luck and Happy reading work!

The assignment must be completed before **11:59 pm on Tuesday, Oct 22, 2024**. Upload your submission as one ZIP file on Canvas.

This assignment consists of two parts: a writing question and a coding question. For the writing question, your submission should be typeset, and all answers must be clearly legible. You are required to submit a PDF file containing your responses to the writing question (you may directly type your answer here and submit the file). For the coding question, you need to submit your code along with the running results, including the reward plot and a testing GIF file. Clearly specify which algorithm you are using to generate these plots. The coding question is divided into two sections: the fundamental coding questions (`hw1_basic.ipynb`) and the advanced coding questions (`hw1_advanced.ipynb`). In general, it is recommended that you complete the writing questions first, followed by the basic coding questions, and finally the advanced coding questions. Note that prior knowledge may be required to complete the advanced section. **Please ensure that you start this assignment as early as possible, as it will take time to run and fine-tune your code.**

The assignment must be completed by each student alone. Collaboration with other students is strictly prohibited. **Questions may be asked on Ed in case of unclear formulations. Do not post answers, partial answers or hints to answers!** Furthermore, we expect all students to adhere to the standards of academic integrity of the Georgia Institute of Technology.

1 Reinforcement Learning Fundamentals (10 points)

1. (3 points) Describe the difference between *supervised learning* and *reinforcement learning*. Use two application examples to highlight the difference.
2. (4 points) Describe how one of your examples from the first question can be formalized with a MDP. State all formal parts of the MDP and sketch how the formal structure relates to your problem.
3. (3 points) When dealing with *infinite horizon problems*, it is common to include a *discount factor*. Describe, in your own words, the role the discount factor plays in the reinforcement learning problem and what problem might occur if the discount factor is set to 1, and what problem might occur if the discount factor is set to 0.

2 Multiple Choices Questions (16 points)

For each of the following multiple-choice questions, select the correct answer(s). There may be more than one correct option per question. To simplify grading, the multiple-choice section will be uploaded to Canvas. **Complete this part on Canvas, as answers not submitted there will not be graded! We will open this section on Canvas as soon as possible!**

1. (2 points) Suppose we are having a well-defined finite grid world MDP problem with terminal state included. For each state, you will have your reward R defined properly and $R \geq 0$ for every grid. Suppose that we want to re-design the reward function there. For which of the following new reward functions would **guarantee** the optimal policy **remain unchanged**? Let $R(s, a, s')$ be the original reward function.
 - ☐ $R_1(s, a, s') = 10R(s, a, s')$
 - ☐ $R_2(s, a, s') = 1 + R(s, a, s')$
 - ☐ $R_3(s, a, s') = R(s, a, s')^2$
 - ☐ $R_4(s, a, s') = 1$
 - ☐ None
2. (2 points) In TD learning with linear VFA (value function approximation), which of the following statement is true?
 - ☐ To present a value function with a weighted linear combination of features, we can write $V(s; w) = x(s)^T w(s)$
 - ☐ For SARSA using VFA, $\Delta w = \alpha(r + \gamma Q(s', \hat{a}'; w) - Q(s, \hat{a}; w)) \nabla_w Q(s', \hat{a}'; w)$
 - ☐ Consider we are using non-linear value function approximation, then during the weights update, we just need to re-calculate $\nabla_w Q(s, a; w)$.
 - ☐ None of the above
3. (2 points) Which statement is true about Policy gradient?
 - ☐ Policy gradient can be computed with automatic differentiation when policy is a neural network. True
 - ☐ Policy gradient is on-policy. True
 - ☐ Policy gradient can also derive off-policy variants using importance sampling. True
 - ☐ In Vanilla policy gradient, the condition number is small when computing the gradient. False
 - ☐ In covariant/natural policy gradient, we re-scale the gradient properly to improve the optimization. True
4. (2 points) Which statement is true about Policy gradient? (Hint: two options are correct.)
 - ☐ Policy gradient will have a low variance since gradient is unbiased.

- ☐ Policy gradient will have a high variance since huge amounts of rollouts will introduce noisy gradients.
 - ☐ Using a larger batch in Policy gradient will lead to high variance.
 - ☐ Tweaking the learning rate such as using adaptive step size rules can reduce the variance in Policy gradient.
 - ☐ Having a fixed learning rate can reduce the variance in Policy gradient.
5. (2 points) Which statement is true about Policy gradient?
- ☐ In Policy gradient, we subtract the baseline to reduce the variance. True
 - ☐ In Policy gradient, subtracting the baseline may introduce bias. False
 - ☐ In policy gradient, The baseline is essentially a proxy for the expected actual return. True
 - ☐ None of the above
6. (2 points) Which statement is true about Actor-critic?
- ☐ Actor-Critic is a Temporal Difference(TD) version of Policy gradient.
 - ☐ Both actor and critic are improving overtime.
 - ☐ Consider actor-critic with discount, we prefer to discount the whole gradient because we care about the whole trajectory. False
 - ☐ We usually design actor and critic with shared features when states are high dimensional such as image state representation. True
 - ☐ In online actor-critic, we have both synchronized and asynchronous version.
7. (2 points) What statements are true about **on-policy** algorithms?
- ☐ They are able to reuse data from different policies in the past directly
 - ☐ They are **not** able to reuse data from different policies in the past directly
 - ☐ They are always more sample efficient than off-policy algorithms
 - ☐ They cannot be used in discrete state-action spaces
 - ☐ None of the above
8. (2 points) Why do we need the log ratio trick in REINFORCE?
- ☐ Because we cannot differentiate the gradient estimator otherwise
 - ☐ Because the original authors did not like automatic differentiation
 - ☐ Because the gradient estimate has too much stochasticity
 - ☐ Because it helps with exploration
 - ☐ None of the above

3 Calculation Questions on Tabular Q-Learning(8 points)

1. (4 points) Consider a system with two states and two actions. You perform actions and observe the rewards and transitions listed below. Each step lists the current state, reward, action and resulting transition as $S_i; R = r; a_k : S_i \rightarrow S_j$. Perform Q-learning using a learning rate $\alpha = 0.5$ and discount factor of $\gamma = 0.5$ for each step. The Q entries are initialized to zero. After step 3, what is $Q(S_2, a_1)$?

Step 1: $S_1; R = -10; a_1 : S_1 \rightarrow S_1$

Step 2: $S_1; R = -10; a_2 : S_1 \rightarrow S_2$

Step 3: $S_2; R = +20; a_1 : S_2 \rightarrow S_1$

2. (4 points) In the same setup as the previous question, what if you are performing SARSA algorithm and you are having an extra step 4 as followed. What is $Q(S_2, a_1)$ after performing the 4 steps?

Consider a system with two states and two actions. You perform actions and observe the rewards and transitions listed below. Each step lists the current state, reward, action and resulting transition as $S_i; R = r; a_k : S_i \rightarrow S_j$. Perform SARSA using a learning rate $\alpha = 0.5$ and discount factor of $\gamma = 1$ for each step. The Q entries are initialized to zero.

Step 1: $S_1; R = -10; a_1 : S_1 \rightarrow S_1$

Step 2: $S_1; R = -10; a_2 : S_1 \rightarrow S_2$

Step 3: $S_2; R = +20; a_1 : S_2 \rightarrow S_1$

Step 4: $S_1; R = +10; a_1 : S_1 \rightarrow S_2$

4 Analyzing a Modern Algorithm: DDPG (14 points)

In this section, we will work our way towards the formulation of a modern high performance reinforcement learning algorithm. As we shall see, this algorithm is ultimately derived from similar principles as the reinforcement learning algorithms discussed in class, with adjustments to increase its stability and efficiency. In particular, we will discuss the Deep Deterministic Policy Gradient (DDPG) algorithm. The DDPG algorithm can be found in a simplified form in Figure 1.

Algorithm 1 Deep Deterministic Policy Gradient

```

1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi$ , empty replay buffer  $\mathcal{D}$ 
2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta$ ,  $\phi_{\text{targ}} \leftarrow \phi$ 
3: repeat
4:   Observe state  $s$  and select action  $a = \text{clip}(\mu_{\theta}(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$ 
5:   Execute  $a$  in the environment
6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
8:   If  $s'$  is terminal, reset environment state.
9:   if it's time to update then
10:    for however many updates do
11:      Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
12:      Compute targets
          
$$y(r, s', d) = r + \gamma(1 - d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$$

13:      Update Q-function by one step of gradient descent using
          
$$\nabla_{\phi} \frac{1}{|B|} \sum_{(s, a, r, s', d) \in B} (Q_{\phi}(s, a) - y(r, s', d))^2$$

14:      Update policy by one step of gradient ascent using
          
$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi}(s, \mu_{\theta}(s))$$

15:      Update target networks with
          
$$\begin{aligned} \phi_{\text{targ}} &\leftarrow \rho \phi_{\text{targ}} + (1 - \rho) \phi \\ \theta_{\text{targ}} &\leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta \end{aligned}$$

16:    end for
17:  end if
18: until convergence

```

Figure 1: DDPG Algorithm

3. (5 points) State and explain one alternate approach to solve the issue raised in the previous part of this question, using pseudo-code to illustrate your idea. What is one benefit and one drawback of this approach compared to DDPG?

4. (2 points) The policy portion of the update is not the same as the standard policy gradient, and instead makes an assumption on some property of the Q function.

State the assumption on Q that DDPG is using here for their policy update, and explain how it justifies the expression in gradient ascent.

5 Reinforcement Learning Application (18 points)

Reinforcement learning method can be widely used in many different field. For this question, we will explore deep reinforcement learning methods in traffic signal control (TSC) applications. Consider a simple TSC problem. In a single four-way intersection, each direction is controlled with green, red and yellow phases by a traffic light. The main goal of TSC is providing safe, effective and reliable transportation systems to participants.

1. (2 points) Please briefly state that why TSC problem can be described as MDP.
2. (4 points) Dealing with the TSC problem on simulators by using RL algorithms requires a good problem formulation in several parts: state, action, reward definitions and neural network structure. In this single four-way intersection situation, please describe what is your state and action in this problem and give a concrete example of each. You need to describe your state and action clearly on what each feature in your state and action represent of and how you are extracting this information.

5. (3 points) Supposed that we are using Policy Optimization method to solve the CMDP problem describing above. What could be a problem compared to MDP?

6. (2 points) Please briefly describe one intuition to solve the problem in the previous question.

6 Coding Questions (65 points)

This section provides instructions for the coding assignment. The coding questions are divided into two parts: basic coding questions (covering algorithms taught in class) and advanced coding questions (implementing a commonly used model-free algorithm, not covered in class but guided step-by-step). Please clone the repository from <https://github.com/pairlab/cs8803drl-fall24.git> to access the code. Follow the instructions provided to complete the coding assignment step by step.

Before beginning the coding tasks, follow the instructions in the `README.md` file to set up a Python environment for running the code. You may use either a GPU or CPU to train your algorithm.

Notice: All hyperparameters inside the Hyper-parameters Tuning cell may require adjustment. Please allocate sufficient time to complete the coding questions.

You are free to add any additional functions or methods to the code if you find that simply completing the TODO sections is not sufficient to make the code run properly.

6.1 Basic Coding Questions (40 points)

For the basic coding question, the code is provided in `hw1_basic.ipynb`. Ensure that you run all cells in this notebook and save the reward plot and GIF file as `{algo_name}_returns.png` and `{algo_name}_policy.gif`, respectively. Follow the instructions for each algorithm to complete the corresponding functions in the `/src` directory.

If you run the code cells in `hw1_basic.ipynb`, `{algo_name}_returns.png` and `{algo_name}_policy.gif` will be automatically generated under the `artifacts/` folder. In the basic coding task, we are using the "LunarLander-v2" environment with discrete action space. You are encouraged to try different discrete or continuous action environments in Gym. You can also explore more complex tasks from Gym-Atari or MuJoCo (note that these may require additional package installations).

1. (10 points) **REINFORCE algorithm:**

Completing the code in `src/network.py` and `src/reinforce.py` based on the TODO instructions should be sufficient to run the REINFORCE algorithm. You can tune the parameters inside the Hyper-parameters Tuning cell in the `hw1_basic.ipynb` file if you find that the provided values are not sufficient to achieve the desired results.

2. (10 points) **Simple Q-iteration (no experience replay + target network):**

Completing the code in `src/network.py` and `src/q_iter.py` based on the TODO instructions should be sufficient to run the Q-iteration algorithm. You can tune the parameters inside the Hyper-parameters Tuning cell in the `hw1_basic.ipynb` file if you find that the provided values are not sufficient to achieve the desired results.

3. (10 points) **DQN (Deep Q-learning + experience replay + target network):**

Completing the code in `src/network.py` and `src/dqn.py` based on the TODO instructions should be sufficient to run the DQN algorithm. You can tune the parameters inside the

Hyper-parameters Tuning cell in the `hw1_basic.ipynb` file if you find that the provided values are not sufficient to achieve the desired results.

4. (10 points) **Actor-Critic:**

Completing the code in `src/network.py` and `src/ac.py` based on the TODO instructions should be sufficient to run the Actor-Critic algorithm. You can tune the parameters inside the **Hyper-parameters Tuning** cell in the `hw1_basic.ipynb` file if you find that the provided values are not sufficient to achieve the desired results.

6.2 Advanced Coding Questions (25 points)

For the advanced coding question, the code is provided in `hw1_advanced.ipynb`. Ensure that you complete and run all cells in this notebook and save the reward plot `sac_returns.png`.

In the advanced coding task, we are using the "Pendulum-v1" environment with continuous action space. You are encouraged to try different continuous action environments in Gym. You can also explore more complex tasks from Gym-Atari or MuJoCo (note that these may require additional package installations).

1. (25 points) **Soft Actor-Critic:**

Complete the code in `hw1_advanced.ipynb` based on the TODO instructions. You can tune the parameters inside the **Hyper-parameters Tuning** cell in the `hw1_advanced.ipynb` file if you find that the provided values are not sufficient to achieve the desired results.

Checkbox for Assignment Submission

Congratulations on completing Assignment 1! We have created a checklist for you to double-check your submission. The files you need to submit are:

- ☐ **hw1.zip:** a zip file with your completed PDF and code.
 - ☐ **Python files:**
`ac.py`, `dqn.py`, `networks.py`, `q_iter.py`, `reinforce.py`
 - ☐ **Jupyter Notebook files:**
`hw1_basic.ipynb`, `hw1_advanced.ipynb`
 - ☐ **Coding results:**
`ac_policy.gif`, `ac_returns.png`,
`dqn_policy.gif`, `dqn_returns.png`,
`q_iteration_policy.gif`, `q_iteration_returns.png`,
`reinforce_policy.gif`, `reinforce_returns.png`,
`sac_returns.png`, `sac_policy.gif` (optional)
- ☐ **Writing assignment PDF file:** `CS8803_DRL_A1.pdf`
- ☐ **hw1_basic.pdf:**
Convert the `hw1_basic.ipynb` to a PDF file **containing your running results.**
- ☐ **hw1_advanced.pdf:**
Convert the `hw1_advanced.ipynb` to a PDF file **containing your running results.**
- ☐ **Completing the Multiple Choices Questions on Canvas.**