

## Assignment 4

**Due Date** November 20, 2023.

**Submission Instructions** Please submit your response in two files via Canvas:

- a pdf file that includes your answers to all the problems, and
- a zip file that includes your completed Jupyter notebooks.

The completed Jupyter notebooks should include missing code as well as the code output.

**Accompanying Code** The assignment includes an accompanying zip file with starter code for the programming problems. To execute and complete this starter code you will need Python 3 and Jupyter. You can use either your local Python 3 run-time or the one provided by [Google Colaboratory](#) (available using your Rice NetID). If you have any questions about working with the code, please contact the teaching staff.

**Collaboration Policy** Collaborating on assignments is permitted, provided the submission lists the students who you collaborated with. As per the [Rice Honor System Handbook](#), this means that students are allowed to develop answers to specific problems together and check answers with each other. However, collaboration does not confer the right for students to submit the exact same document—students must write down the answers themselves. While the “core” of the answer can be the same, the wording cannot be identical unless precise wording is necessary to answer the question. Students must be able to demonstrate that they worked together when developing responses and that one student did not copy off the other. This means that students are barred from dividing questions among themselves.

**Citation Policy** Students don’t have an obligation to cite class slides, lectures, or class textbooks on assignments; these are considered common knowledge. An academic citation style is required for all other sources (such as research articles and blogs).

**Late Policy** Assignments should be turned on time via Canvas. Assignments handed in late will be marked off 10% per day. Assignments more than 3 days late will not be accepted. In turn, you can expect the teaching staff to grade your assignments and provide feedback in a timely manner.

**Grading** The assignment is worth 20% of your final grade.

**Updates to the Assignment** In case there are any updates to the assignment (e.g., additional clarifications, typo fixes, hints, etc.), they will be indicated via the following table.

Version	Date	Note
v1	November 5	Assignment released

## 1. Deep Deterministic Policy Gradient

Starter code for this problem is provided in the accompanying **Jupyter** notebook. For each part, include the resulting plots with 2-4 sentence summary of these plots in your **pdf** submission.

### 1.1 30pts

Implement the general recipe for Deep Deterministic Policy Gradient (DDPG) algorithm with function approximation. Use this recipe to learn the optimal policy for the [Lunar Lander \(v2\)](#) environment from [Gymnasium](#).

### 1.2 30pts

Report results of five runs of your algorithm and summarize the performance in a single plot. Discuss the root cause of variance observed in the results and methods to reduce the variance.

### 1.3 15pts

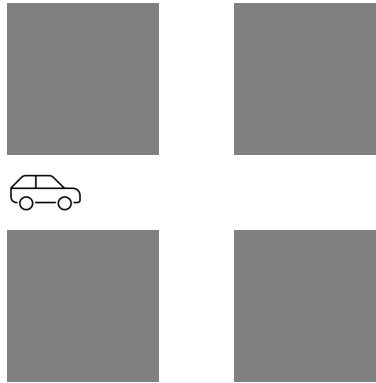
Select any one hyperparameter of the algorithm and study its effect on the agent's performance. You should try at least five different values of this hyperparameter and report on the observed performance.

### 1.4 10pts

While applying RL in real world, you may consider the use of off-the-shelf implementation of RL algorithms. The package [Stable Baselines 3](#) aims to provide a set of reliable implementations of RL algorithms in PyTorch.

Familiarize yourself with this package and use its implementation of [DDPG](#) to learn the optimal policy for the [Lunar Lander \(v2\)](#) environment.

## 2. Principle of Maximum Entropy

**2.1****5pts**

As we have seen in the class, entropy is often used as a regularizer in reinforcement and imitation learning algorithms. Provide the mathematical definition of entropy  $H(p)$  of a distribution  $p(x)$ .

**2.2****10pts**

Consider a car approaching a four-way intersection, as shown above. The car's driver is equally likely to go straight or make a turn, i.e.,

$$\Pr(\text{turn left}) + \Pr(\text{turn right}) = \Pr(\text{go straight}), \text{ and} \quad (1)$$

will not make a U-turn. Given this information and using the principle of maximum entropy, analytically derive the probability distribution over the decision taken by the car driver at the four-way intersection. Does the maximum entropy distribution match your intuition?