Assignment 1: Imitation Learning

Due September 12, 11:59 pm

The goal of this assignment is to experiment with imitation learning, including direct behavior cloning and the DAgger algorithm. In lieu of a human demonstrator, demonstrations will be provided via an expert policy that we have trained for you. Your goals will be to set up behavior cloning and DAgger, and compare their performance on a few different continuous control tasks from the OpenAI Gym benchmark suite. Turn in your report and code as described in Section 4.

The starter-code for this assignment can be found at

```
https://github.com/berkeleydeeprlcourse/homework_fall2022
```

You have the option of running the code either on Google Colab or on your own machine. Please refer to the README for more information on setup.

1 Behavioral Cloning

1. The starter code provides an expert policy for each of the MuJoCo tasks in OpenAI Gym. Fill in the blanks in the code marked with Todo to implement behavioral cloning. A command for running behavioral cloning is given in the Readme file.

We recommend that you read the files in the following order. For some files, you will need to fill in blanks, labeled TODO.

- scripts/run_hw1.py (read-only)
- infrastructure/rl_trainer.py
- agents/bc_agent.py (another read-only file)
- policies/MLP_policy.py
- infrastructure/replay_buffer.py
- infrastructure/utils.py
- infrastructure/pytorch_utils.py
- 2. Run behavioral cloning (BC) and report results on two tasks: the Ant environment, where a behavioral cloning agent should achieve at least 30% of the performance of the expert, and one environment of your choosing where it does not. Here is how you can run the Ant task:

```
python cs285/scripts/run_hw1.py \
--expert_policy_file cs285/policies/experts/Ant.pkl \
--env_name Ant-v4 --exp_name bc_ant --n_iter 1 \
--expert_data cs285/expert_data/expert_data_Ant-v4.pkl \
--video_log_freq -1
```

When providing results, report the mean and standard deviation of your policy's return over multiple rollouts in a table, and state which task was used. When comparing one that is working versus one that is not working, be sure to set up a fair comparison in terms of network size, amount of data, and number of training iterations. Provide these details (and any others you feel are appropriate) in the table caption.

Note: What "report the mean and standard deviation" means is that your eval_batch_size should be greater than ep_len, such that you're collecting multiple rollouts when evaluating the performance of your trained policy. For example, if ep_len is 1000 and eval_batch_size is 5000, then you'll be collecting approximately 5 trajectories (maybe more if any of them terminate early), and the logged

Eval_AverageReturn and Eval_StdReturn represents the mean/std of your policy over these 5 rollouts. Make sure you include these parameters in the table caption as well.

Tip: To generate videos of the policy, remove the flag --video_log_freq -1 However, this is slower, and so you probably want to keep this flag on while debugging.

3. Experiment with one set of hyperparameters that affects the performance of the behavioral cloning agent, such as the amount of training steps, the amount of expert data provided, or something that you come up with yourself. For one of the tasks used in the previous question, show a graph of how the BC agent's performance varies with the value of this hyperparameter. In the caption for the graph, state the hyperparameter and a brief rationale for why you chose it.

2 DAgger

1. Once you've filled in all of the TODO commands, you should be able to run DAgger.

```
python cs285/scripts/run_hw1.py \
--expert_policy_file cs285/policies/experts/Ant.pkl \
--env_name Ant-v4 --exp_name dagger_ant --n_iter 10 \
--do_dagger --expert_data cs285/expert_data/expert_data_Ant-v4.pkl \
--video_log_freq -1
```

2. Run DAgger and report results on the two tasks you tested previously with behavioral cloning (i.e., Ant + another environment). Report your results in the form of a learning curve, plotting the number of DAgger iterations vs. the policy's mean return, with error bars to show the standard deviation. Include the performance of the expert policy and the behavioral cloning agent on the same plot (as horizontal lines that go across the plot). In the caption, state which task you used, and any details regarding network architecture, amount of data, etc. (as in the previous section).

3 Turning it in

1. **Submitting the PDF.** Make a PDF report containing: Table 1 for a table of results from Question 1.2, Figure 1 for Question 1.3. and Figure 2 with results from question 2.2.

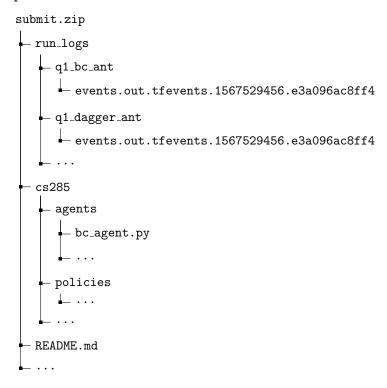
You do not need to write anything else in the report, just include the figures with captions as described in each question above. See the handout at

```
http://rail.eecs.berkeley.edu/deeprlcourse/static/misc/viz.pdf
```

for notes on how to generate plots.

- 2. Submitting the code and experiment runs. In order to turn in your code and experiment logs, create a folder that contains the following:
 - A folder named run_logs with experiments for both the behavioral cloning (part 2, not part 3) exercise and the DAgger exercise. Note that you can include multiple runs per exercise if you'd like, but you must include at least one run (of any task/environment) per exercise. These folders can be copied directly from the cs285/data folder into this new folder. Important: Disable video logging for the runs that you submit, otherwise the files size will be too large! You can do this by setting the flag --video_log_freq -1
 - The cs285 folder with all the .py files, with the same names and directory structure as the original homework repository. Also include the commands (with clear hyperparameters) that we need in order to run the code and produce the numbers that are in your figures/tables (e.g. run "python run_hw1_behavior_cloning.py -ep_len 200" to generate the numbers for Section 2 Question 2) in the form of a README file.

As an example, the unzipped version of your submission should result in the following file structure. Make sure that the submit.zip file is below 15MB and that they include the prefix $q1_{-}$ and $q2_{-}$.



- 3. If you are a Mac user, do not use the default "Compress" option to create the zip. It creates artifacts that the autograder does not like. You may use zip -vr submit.zip submit -x "*.DS_Store" from your terminal.
- 4. Turn in your assignment on Gradescope. Upload the zip file with your code and log files to $\mathbf{HW1}$ \mathbf{Code} , and upload the PDF of your report to $\mathbf{HW1}$.