

UNIVERSITY OF VIRGINIA
CS 6501-012: LEARNING IN ROBOTICS
HOMEWORK 4
DUE: 05/07/25 TUE BY 11.59 PM

Instructions

Read the following instructions carefully before beginning to work on the homework.

- All the files required for this homework (starter code, reference PDFs etc.) are available in the following directory: **Canvas » Files » hw4**.
- You will submit solutions typeset in \LaTeX on Gradescope. You can use `hw_template.tex` from the course website (under assignments).
- Please start a new problem on a fresh page and mark all the pages corresponding to each problem. Failure to do so may result in your work not graded completely.
- Clearly indicate the name and UVA email ID of all your collaborators on your submitted solutions.
- **For each problem in the homework, you should mention the total amount of time you spent on it. This helps us gauge the perceived difficulty of the problems.**
- You can be informal while typesetting the solutions, e.g., if you want to draw a picture feel free to draw it on paper clearly, click a picture and include it in your solution. Do not spend undue time on typesetting solutions.
- You will see an entry of the form “HW 4 PDF” where you will upload the PDF of your solutions. You will also see entries like “HW 4 Problem X Code” where you will upload your solution for the respective problems. **For each programming problem, you should create a fresh Python file.** This file should contain all the code to reproduce the results of the problem and you will upload the `.py` file to Gradescope. If we have installed Autograder for a particular problem, you will use the Autograder. Name your file to be the same filename as stated in the respective problem statement.
- **You should include all the relevant plots in the PDF, without doing so you will not get full credit.** You can, for instance, export your Jupyter notebook as a PDF (you can also use text cells to write your solutions) and export the same notebook as a Python file to upload your code.

- **Your PDF solutions should be completely self-contained. We will run the Python file to check if your solution reproduces the results in the PDF.**

Credit The points for the problems add up to 100. You only need to solve for 100 points to get full credit, i.e., your final score will be $\min(\text{your total points}, 100)$.

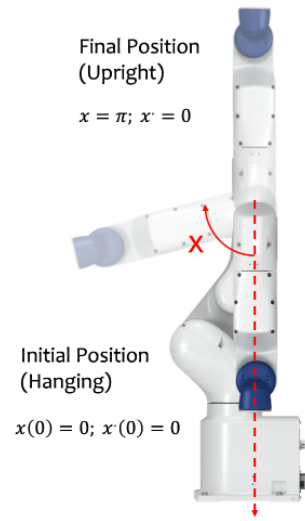
1 **Problem 1 (Robotic Arm Stabilization, 50 points).** You will develop and train a
 2 controller to balance a robotic arm, modeled similarly to a damped pendulum, to
 3 maintain an upright position. The arm has a single joint, akin to a human elbow,
 4 and must be controlled to stabilize against gravity and other disturbances.

5 Rather than using a complex robotic arm simulator, we will employ the following
 6 simplified dynamics to simulate the arm's behavior (think of the equation below as
 7 our simulator for this problem):

$$ml^2\ddot{x} + b\dot{x} + mgl \sin(x) = u$$

8 where:

- x is the angle of the arm relative to the hanging downward position (in radians), measured from the downward vertical line.
- \ddot{x} and \dot{x} are the angular acceleration and velocity,
- u is the torque applied at the joint,
- $g = 9.8 \text{ m/s}^2$ (acceleration due to gravity),
- $m = 1 \text{ kg}$ (mass of the arm),
- $l = 1 \text{ meter}$ (length of the arm from the joint to the center of mass),
- $b = 0.1 \text{ Nm/s}$ (damping coefficient representing joint friction).



10
 11 The robotic arm starts from a downwards hanging position with $x(0) = \dot{x}(0) = 0$
 12 and the task is to move it to and maintain it in the upright position where $x = \pi$ and
 13 $\dot{x} = 0$. The control input u is constrained such that:

$$|u| \leq 1.$$

14 Implement this controller using policy gradients and a neural network. The reward
 15 at a state x and control u is

$$r(x, \dot{x}, u) = -\frac{1}{2} \left[(\pi - x)^2 + \dot{x}^2 + \frac{1}{100} u^2 \right].$$

16 We have provided you some example code on Canvas (p1.py); feel free to modify
 17 this code as you wish. But read the comments inside the code carefully before
 18 beginning to write your solution.

- 19 (a) **(10 points)** How is the stochastic controller $u_\theta(\cdot | x)$ implemented in
 20 the code ? What is the probability distribution, how is the log-likelihood
 21 $\log u_\theta(u|x)$ computed ? How is the constraint $|u| \leq 1$ enforced in the
 22 implementation ?
 23 (b) **(30 points)** Implement code to train the policy using policy gradient. In the
 24 PDF, explain in detail:

- 1 • [5] How many trajectories you collected per iteration and the total
- 2 number of iterations you trained for. Include code snippets showing
- 3 the sampling loop.
- 4 • [10] How did you compute the policy gradient ? Include a code
- 5 snippet that clearly shows and breaks down how you computed the
- 6 log-likelihoods, multiplied by the returns.
- 7 • [10] What form of baseline did you use and how did you implement it
- 8 ?
- 9 • [5] How gradient ascent was performed ? What optimizer did you use?
- 10 What learning rate? What did you try, and what worked ?
- 11 (c) (10 points) First report and plot the cumulative reward over 1000 time-steps
- 12 as a function of parameter updates to θ for [5] points. Then Modify your
- 13 code to change the mass to $m = 2$ and evaluate the trained policy (with
- 14 $m = 1$ on this new dynamics. Report the cumulative reward and comment
- 15 on generalization. Only the evaluation phase uses $m = 2$. You should
- 16 freeze the weights for $m = 1$ and do not retrain. [5]

17 **Problem 2 (Q-Learning, 50 points).** You will write code for Q-learning with
 18 the DDQN trick in this problem for a simple environment called the CartPole
 19 <https://stanford.edu/jeffjar/cartpole.html>. We have provided some example code
 20 that uses PyTorch for Q-learning. You need to fill in the functions for epsilon-greedy
 21 exploration and the optimization objective. The code is given at p2.py on Canvas;
 22 feel free to modify this code as you wish.

23 We will be using OpenAI Gym's version of CartPole, read the code at
 24 https://github.com/openai/gym/blob/master/gym/envs/classic_control/cartpole.py to
 25 understand the environment. The class for the q -function in the example code looks
 26 as follows. The neural network s.m is a two-layer neural network with ReLU non-
 27 linearity. Notice that we have structured the network not as $q(x, u) : X \times U \rightarrow \mathbb{R}$
 28 but instead as $q(x) : X \rightarrow U$. This way the neural network returns the q -value
 29 for all controls u with a single call to `q_t.forward`. You should implement the
 30 epsilon-greedy strategy to pick a control in the function `q_t.control`.

```

31
32
33 class q_t(nn.Module):
34     def __init__(s, xdim, udim, hdim=16):
35         super().__init__()
36         s.xdim, s.udim = xdim, udim
37         s.m = nn.Sequential(
38             nn.Linear(xdim, hdim),
39             nn.ReLU(True),
40             nn.Linear(hdim, udim),
41         )
42     def forward(s, x):
43         return s.m(x)
44
45     def control(s, x, eps=0):

```

```

1      # 1. get q values for all controls
2      q = s.m(x)
3
4      # eps-greedy strategy to choose control input
5      # note that for eps=0
6      # you should return the correct control u
7      return u
8

```

9 Read the rollout function carefully. It takes a q -network and runs it for T timesteps to return a trajectory. You should add this trajectory to the replay buffer.

```

11
12 def rollout(e, q, eps=0, T=1000):
13     traj = []
14
15     x = e.reset()
16     for t in range(T):
17         u = q.control(th.from_numpy(x).float().unsqueeze(0),
18                     eps=eps)
19         u = u.int().numpy().squeeze()
20
21         xp, r, d, info = e.step(u)
22         t = dict(x=x, xp=xp, r=r, u=u, d=d, info=info)
23         x = xp
24         traj.append(t)
25         if d:
26             break
27     return traj
28

```

29 You will code up the Bellman error minimization objective with the Double-
30 Q network trick. Hint: You can use the following code to create a copy of the
31 q -network. You can also modify the class `q_t` to create a copy of `s.m` inside it
32 directly.

```

33
34 import copy
35 qc = copy.deepcopy(q)

```

37 Read the main function. We will create an environment `e` using the OpenAI Gym
38 library, then initialize the q -network and create an optimizer (in this case Adam) to
39 update the parameters of the q -network. The power of PyTorch lies in being able to
40 call `f.backward()` to compute the gradient of whatever objective that depends on the
41 parameters of the q -function. The call `optim.step()` updates the parameters of the
42 value-function.

```

43
44 if __name__ == '__main__':
45     e = gym.make('CartPole-v0')
46
47     xdim, udim = e.observation_space.shape[0], \
48                 e.action_space.n
49
50     q = q_t(xdim, udim, 8)
51     optim = th.optim.Adam(q.parameters(), lr=1e-3,

```

```

1         weight_decay=1e-4)
2
3     # this is the replay buffer
4     ds = []
5
6     # collect few random trajectories with
7     # eps=1
8     for i in range(1000):
9         ds.append(rollout(e, q, eps=1, T=200))
10
11    for i in range(1000):
12        q.train()
13        t = rollout(e, q)
14        ds.append(t)
15
16        # perform sgd updates on the q network
17        # need to call zero grad on q function
18        # to clear the gradient buffer
19        q.zero_grad()
20        f = loss(q, ds)
21        f.backward()
22        optim.step()
23    print('Log data to plot')

```

25 During training you should keep track of the average return of the network.
26 You should also fill in this function that evaluates the learnt q -function on a new
27 environment.

```

28
29 def evaluate(q):
30     # 1. create a new environment e
31     # 2. run the learned q network for 100 trajectories on
32     # this new environment and report the average discounted
33     # return of these 100 trajectories
34     return r

```

- 36 **(30 points)** Correctly code up all the methods above. Plot the following
- 37 1. **(30 points)** Explain your implementation with code snippets.
 - 38 2. **(10 points)** Every 1000 gradient steps, evaluate the current q -network by
39 running the learned policy for 10 episodes on the **training** environment. Plot
40 the average total return across these episodes over the course of training.
 - 41 3. **(10 points)** Similarly, evaluate the learned policy on a separate evaluation
42 environment every 1000 gradient steps. Plot the average total return across
43 10 episodes over training.
- 44 The maximum achievable average reward for this environment is 200. You can call
45 the `render()` function of the environment to see your trained policy on the Cartpole
46 in action.