

Exercise 6

Tommaso Brumani - 100481325
ELEC-E8125 - Reinforcement Learning

November 6, 2022

Task 1

The actor-critic algorithm was implemented for the Inverted Pendulum environment by modifying the **pg_ac.py** file as follows:

- In the **Policy.__init()** function, the following line was added:

```
1  # Implement actor_logstd as a learnable parameter
2  # Use log of std to make sure std doesn't become negative
   during training
3  self.actor_logstd = nn.parameter.Parameter(torch.tensor(np.
   zeros(action_dim), device=device))
```

- In the **Policy.forward()** function, the following line was added:

```
1  # Create a Normal distribution with mean of 'action_mean'
   and standard deviation of 'action_logstd', and return
   the distribution
2  probs = Normal(action_mean, action_std)
```

- In the `PG.update()` function, the following lines were added:

```
1 ##### Your code starts here. #####
2
3 # calculate values for current and next states
4 states_values = self.value.forward(states)
5 next_states_values = self.value.forward(next_states)
6
7 # set next state values to zero when terminal
8 dones_mask = dones == False
9 next_states_values *= dones_mask
10
11 # calculate TD target
12 td_target = rewards + self.gamma * next_states_values
13
14 # calculate critic loss
15 critic_loss = torch.mean(torch.square(states_values -
    td_target.detach()))
16
17 # calculate advantage
18 adv = td_target - states_values
19
20 # normalize advantage
21 (std, mean) = torch.std_mean(adv)
22 norm_adv = (adv - mean)/std
23
24 # compute policy loss for batch
25 actor_loss = - torch.mean(action_probs.squeeze(-1) *
    norm_adv.detach())
26
27 # compute combined loss
28 comb_loss = actor_loss + critic_loss
29
30 # backpropagate gradients
31 comb_loss.backward()
32 self.optimizer.step()
33
34 # empty optimizer gradients
35 self.optimizer.zero_grad()
36
37 ##### Your code ends here. #####
```

- In the `PG.get_action()` function, the following lines were added:

```
1 ##### Your code starts here. #####
2
3 # get action distribution following policy for give
  observation
4 action_distribution = self.policy.forward(x)
5
6 # select action
7 if evaluation:
8     # mean if testing
9     action = action_distribution.mean
10 else:
11     # random sample otherwise
12     action = action_distribution.sample()
13
14 # get log probability of chosen action
15 act_logprob = action_distribution.log_prob(action)
16
17 ##### Your code ends here. #####
```

As a result, the performance plot reported in Figure 1 was obtained.

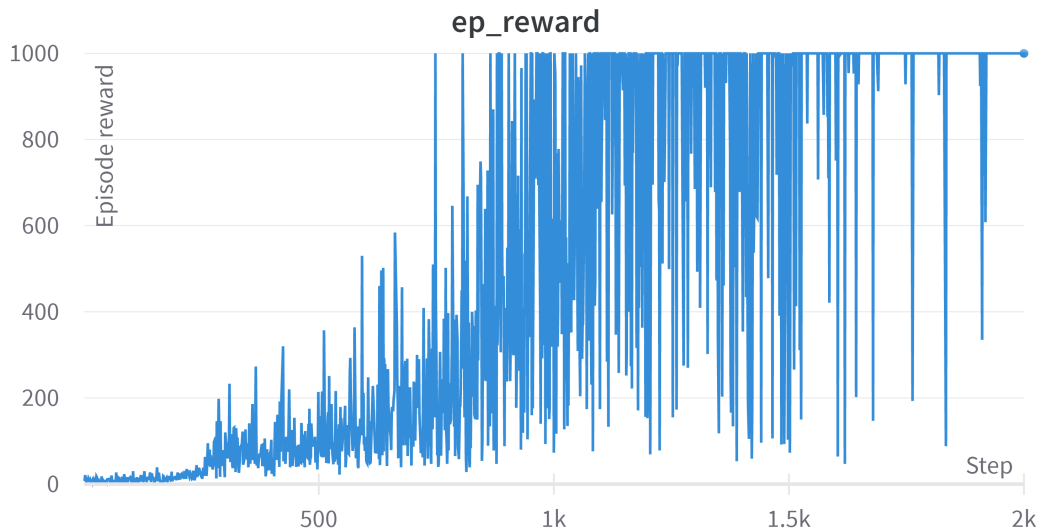


Figure 1: Performance of the PG actor-critic model.

Question 1.1

In REINFORCE with baseline, a good choice for the baseline was an estimation of the state value function $V(s_t)$, which could be learned from the data. In actor-critic methods this becomes the critic model, which is used to determine the TD target.

Question 1.2

The value of Advantage can be intuitively interpreted as the amount by which having chosen action a increases the cumulative return, with respect to the average obtainable from the current state.

In other words, it indicates how optimal the chosen action is with respect to the others.

Question 1.3

Since REINFORCE is a MC approach and the implemented actor-critic model is a TD approach, we can expect the actor-critic to be characterized by an increased bias and reduced variance with respect to the REINFORCE algorithm because of bootstrapping.

Question 1.4

Implementing a form of TD- λ could be used to control the bias/variance tradeoff by setting using values of λ ($\lambda = 0$ corresponding to the current actor-critic approach, and $\lambda = 1$ approximating the MC approach).

Task 2

The deep deterministic policy gradient algorithm was implemented for the Half Cheetah environment by modifying the **ddpg.py** file as follows:

- In the **DDPG.update()** function, the following lines were added:

```
1  ##### Your code starts here. #####
2
3  # compute Q target
4  a_next = self.pi_target(batch.next_state)
5  q_next = self.q_target(batch.next_state, a_next)
6  q_tar = batch.reward + self.gamma * q_next * batch.not_done
7
8  # compute critic loss
9  q_policy = self.q(batch.state, batch.action)
10 critic_loss = torch.mean(torch.square(q_tar.detach() -
    q_policy))
11
12 # update critic's parameter
13 critic_loss.backward()
14 self.q_optim.step()
15
16 # compute actor loss
17 a_policy_new = self.pi(batch.state)
18 q_policy_new = self.q(batch.state, a_policy_new)
19 actor_loss = - torch.mean(q_policy_new)
20
21 # update actor's parameters
22 actor_loss.backward()
23 self.pi_optim.step()
24
25 # update target networks
26 h.soft_update_params(self.q, self.q_target, self.tau)
27 h.soft_update_params(self.pi, self.pi_target, self.tau)
28
29 # empty optimizer gradients
30 self.q_optim.zero_grad()
31 self.pi_optim.zero_grad()
32
33 ##### Your code ends here. #####
```

- In the `DDPG.get_action()` function, the following lines were added:

```
1 ##### Your code starts here. #####
2
3 # caculate action to execute
4 action = self.pi(x)
5
6 # if not evaluation, add gaussain noise
7 if not evaluation:
8     action += expl_noise * torch.randn(action.shape, device
9                                         =device)
10
11 # get correct action shape
12 action = action.squeeze(0)
13 ##### Your code ends here. #####
```

As a result, the performance plot reported in Figure 2 was obtained.

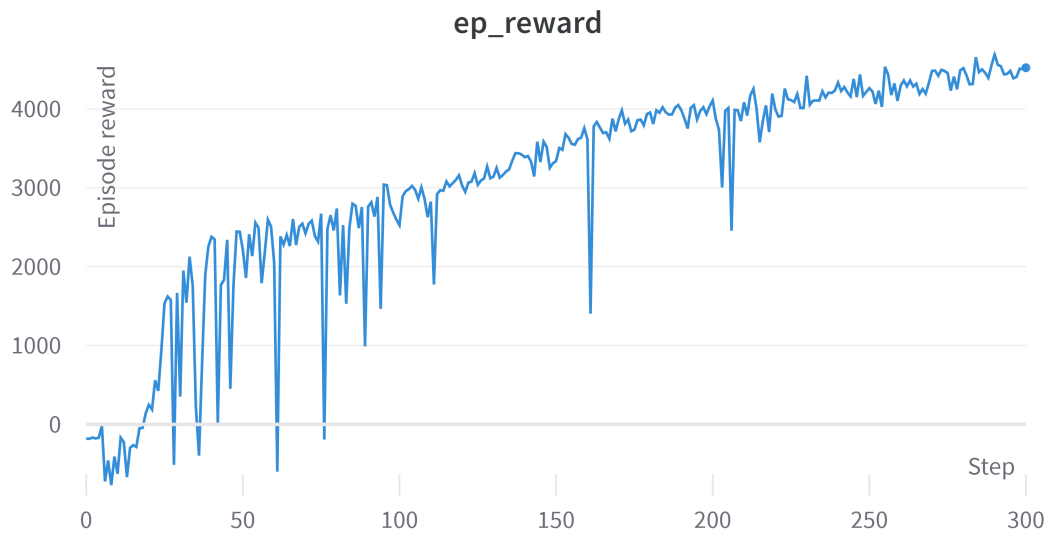


Figure 2: Performance of the DDPG model.

Question 2.1

In DDPG we can use off-policy data because instead of using a whole episode's discounted reward we are able to use the data from single timesteps thanks to bootstrapping. This is made possible by the fact that we use the data sampled from the memory buffer to train the policy network, which in turn we use to determine the estimated optimal Q when calculating the loss functions (similarly to how in Q -learning we used a maximum over the different actions).

Question 2.2

Compared to the actor-critic method of Task 1, which uses the actor to output the mean and variance of a gaussian distribution from which actions are sampled, the DDPG algorithm's actor outputs directly a deterministic value.

While this reduces computational burden in those cases where actions are deterministic, it fails to capture those instances where actions are stochastic, and thus has reduced expressive power.

This also means that the exploration during training is dictated by previously set hyperparameters of the noise distribution, rather than on a probability distribution whose parameters are learned and updated at runtime based on data.