

⌕ B I <> 🔗 🖼️ 📄 📋 📌 🔍 😊 ☰

ECEN 743: Reinforcement Learning

Deep Q-Learning

Code tested using

1. gymnasium 0.27.1
2. box2d-py 2.3.5
3. pytorch 2.0.0
4. Python 3.9.12

1 & 2 can be installed using pip install gymnasium[box2d]

General Instructions

1. This code consists of TODO blocks, read them carefully and complete the blocks

2. Type your code between the following lines

```
##### TYPE YOUR CODE HERE #####
#####
```

3. The default hyperparameters should be able to solve LunarLander-v2

4. You do not need to modify the rest of the code for this assignment, free to do so if needed.

I HAVE USED CARTRPOLE FROM THE GYM LIBRARY

CARTPOLE GYM DOCS: [https://www.gymnasium.dev/environments/classic\\_control/cart\\_pole/](https://www.gymnasium.dev/environments/classic_control/cart_pole/)

ECEN 743: Reinforcement Learning Deep Q-Learning Code tested using 1. gymnasium 0.27.1 2. box2d-py 2.3.5 3. pytorch 2.0.0 4. Python 3.9.12 1 & 2 can be installed using pip install gymnasium[box2d] General Instructions

1. This code consists of TODO blocks, read them carefully and complete each of the blocks

2. Type your code between the following lines

```
##### TYPE YOUR CODE HERE #####
#####
```

3. The default hyperparameters should be able to solve LunarLander-v2

4. You do not need to modify the rest of the code for this assignment, feel free to do so if needed.

I HAVE USED CARTRPOLE FROM THE GYM LIBRARY CARTPOLE GYM DOCS:

[https://www.gymnasium.dev/environments/classic\\_control/cart\\_pole/](https://www.gymnasium.dev/environments/classic_control/cart_pole/)

```
1 !pip install gymnasium[box2d]
2 !pip3 install matplotlib
```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>

Collecting gymnasium[box2d]

Downloading gymnasium-0.28.1-py3-none-any.whl (925 kB)

925.5/925.5 KB 11.6 MB/s eta 0:00:00

Collecting farama-notifications>=0.0.1

Downloading Farama\_Notifications-0.0.4-py3-none-any.whl (2.5 kB)

Requirement already satisfied: importlib-metadata>=4.8.0 in /usr/local/lib/python3.9/dist-packages (from gymnasium[box2d]) (6.1.0)

Requirement already satisfied: typing-extensions>=4.3.0 in /usr/local/lib/python3.9/dist-packages (from gymnasium[box2d]) (4.5.0)

Collecting jax-jumpy>=1.0.0

Downloading jax\_jumpy-1.0.0-py3-none-any.whl (20 kB)

Requirement already satisfied: cloudpickle>=1.2.0 in /usr/local/lib/python3.9/dist-packages (from gymnasium[box2d]) (2.2.1)

Requirement already satisfied: numpy>=1.21.0 in /usr/local/lib/python3.9/dist-packages (from gymnasium[box2d]) (1.22.4)

Collecting swig==4.\*

Downloading swig-4.1.1-py2.py3-none-manylinux\_2\_5\_x86\_64.manylinux1\_x86\_64.whl (1.8 MB)

1.8/1.8 MB 21.4 MB/s eta 0:00:00

Collecting box2d-py==2.3.5

Downloading box2d-py-2.3.5.tar.gz (374 kB)

374.4/374.4 KB 9.5 MB/s eta 0:00:00

Preparing metadata (setup.py) ... done

Collecting pygame==2.1.3

Downloading pygame-2.1.3-cp39-cp39-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl (13.7 MB)

13.7/13.7 MB 21.6 MB/s eta 0:00:00

Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.9/dist-packages (from importlib-metadata>=4.8.0->gymnasium[box2d])

Building wheels for collected packages: box2d-py

**error: subprocess-exited-with-error**

× python setup.py bdist\_wheel did not run successfully.

exit code: 1

→ See above for output.

**note:** This error originates from a subprocess, and is likely not a problem with pip.

Building wheel for box2d-py (setup.py) ... error

**ERROR: Failed building wheel for box2d-py**

Running setup.py clean for box2d-py

Failed to build box2d-py

Installing collected packages: swig, farama-notifications, box2d-py, pygame, jax-jumpy, gymnasium

Running setup.py install for box2d-py ... done

**DEPRECATION: box2d-py was installed using the legacy 'setup.py install' method, because a wheel could not be built for it. A possi**

Attempting uninstall: pygame

Found existing installation: pygame 2.3.0

Uninstalling pygame-2.3.0:

Successfully uninstalled pygame-2.3.0

Successfully installed box2d-py-2.3.5 farama-notifications-0.0.4 gymnasium-0.28.1 jax-jumpy-1.0.0 pygame-2.1.3 swig-4.1.1

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>

Requirement already satisfied: matplotlib in /usr/local/lib/python3.9/dist-packages (3.7.1)

Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (0.11.0)

Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (3.0.9)

Requirement already satisfied: importlib-resources>=3.2.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (5.12.0)

Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (23.0)

Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (2.8.2)

```

Requirement already satisfied: contourpy==1.0.1 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (1.0.1)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (8.4.0)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (4.39.3)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (1.4.4)
Requirement already satisfied: numpy>=1.20 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (1.22.4)
Requirement already satisfied: zipp>=3.1.0 in /usr/local/lib/python3.9/dist-packages (from importlib-resources>=3.2.0->matplotlib) (

```

```

1 import gymnasium as gym
2 import random
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 import torch.optim as optim
7 import argparse
8 import numpy as np
9 from collections import deque, namedtuple
10 import matplotlib.pyplot as plt
11 import base64, io
12
13 # For visualization
14 from gym.wrappers.monitoring import video_recorder
15 from IPython.display import HTML
16 from IPython import display
17 import glob

```

```

1 class ExperienceReplay:
2     """
3     Based on the Replay Buffer implementation of TD3
4     Reference: https://github.com/sfujim/TD3/blob/master/utils.py
5     """
6     def __init__(self, state_dim, action_dim, max_size, batch_size, gpu_index=0):
7         self.max_size = max_size
8         self.ptr = 0
9         self.size = 0
10        self.state = np.zeros((max_size, state_dim))
11        self.action = np.zeros((max_size, action_dim))
12        self.next_state = np.zeros((max_size, state_dim))
13        self.reward = np.zeros((max_size, 1))
14        self.done = np.zeros((max_size, 1))
15        self.batch_size = batch_size
16        self.device = torch.device('cuda', index=gpu_index) if torch.cuda.is_available() else torch.device('cpu')
17
18
19    def add(self, state, action, reward, next_state, done):
20        self.state[self.ptr] = state
21        self.action[self.ptr] = action
22        self.next_state[self.ptr] = next_state
23        self.reward[self.ptr] = reward
24        self.done[self.ptr] = done
25        self.ptr = (self.ptr + 1) % self.max_size
26        self.size = min(self.size + 1, self.max_size)
27
28    def sample(self):
29        ind = np.random.randint(0, self.size, size=self.batch_size)
30
31        return (
32            torch.FloatTensor(self.state[ind]).to(self.device),
33            torch.FloatTensor(self.action[ind]).long().to(self.device),
34            torch.FloatTensor(self.reward[ind]).to(self.device),
35            torch.FloatTensor(self.next_state[ind]).to(self.device),
36            torch.FloatTensor(self.done[ind]).to(self.device)
37        )
38
39

```

```

1 class QNetwork(nn.Module):
2     """
3     Q Network: designed to take state as input and give out Q values of actions as output
4     """
5
6     def __init__(self, state_dim, action_dim):
7         """
8         state_dim (int): state dimension
9         action_dim (int): action dimension
10        """

```

```

11     super(QNetwork, self).__init__()
12     self.l1 = nn.Linear(state_dim, 64)
13     self.l2 = nn.Linear(64, 64)
14     self.l3 = nn.Linear(64, action_dim)
15
16     def forward(self, state):
17         q = F.relu(self.l1(state))
18         q = F.relu(self.l2(q))
19         return self.l3(q)

1 class DQNAgent():
2
3     def __init__(self,
4         state_dim,
5         action_dim,
6         discount=0.99,
7         tau=1e-3,
8         lr=5e-4,
9         update_freq=4,
10        max_size=int(1e5),
11        batch_size=64,
12        gpu_index=0
13    ):
14        """
15        state_size (int): dimension of each state
16        action_size (int): dimension of each action
17        discount (float): discount factor
18        tau (float): used to update q-target
19        lr (float): learning rate
20        update_freq (int): update frequency of target network
21        max_size (int): experience replay buffer size
22        batch_size (int): training batch size
23        gpu_index (int): GPU used for training
24        """
25        self.state_dim = state_dim
26        self.action_dim = action_dim
27        self.discount = discount
28        self.tau = tau
29        self.lr = lr
30        self.update_freq = update_freq
31        self.max_size = max_size
32        self.batch_size = batch_size
33        self.device = torch.device('cuda', index=gpu_index) if torch.cuda.is_available() else torch.device('cpu')
34
35
36        # Setting up the NNS
37        self.Q = QNetwork(state_dim, action_dim).to(self.device)
38        self.Q_target = QNetwork(state_dim, action_dim).to(self.device)
39        self.optimizer = optim.Adam(self.Q.parameters(), lr=self.lr)
40
41        # Experience Replay Buffer
42        self.memory = ExperienceReplay(state_dim,1,max_size,self.batch_size,gpu_index)
43
44        self.t_train = 0
45
46    def step(self, state, action, reward, next_state, done):
47        """
48        1. Adds (s,a,r,s') to the experience replay buffer, and updates the networks
49        2. Learns when the experience replay buffer has enough samples
50        3. Updates target network
51        """
52        self.memory.add(state, action, reward, next_state, done)
53        self.t_train += 1
54
55        # Experience Replay
56        if self.memory.size > self.batch_size:
57            experiences = self.memory.sample()
58            self.learn(experiences, self.discount) #To be implemented
59
60        # Target Network
61        if (self.t_train % self.update_freq) == 0:
62            self.target_update(self.Q, self.Q_target, self.tau) #To be implemented
63
64    def select_action(self, state, epsilon = 0.):
65        if np.random.random() > epsilon:
66            state = torch.tensor(state, dtype=torch.float32, device=self.device) # converting our state to pytorch tensor
67            #self.Q(state) # writing updated tensor to device

```

```

68         self.Q.eval()
69         with torch.no_grad():
70             actions=self.Q.forward(state)
71             self.Q.train()
72             action = torch.argmax(actions).item() # .item() is to convert from tensors to integers
73         else:
74             action=np.random.choice(self.action_dim)
75
76         return action
77
78
79 def learn(self, experiences, discount):
80     """
81     TODO: Complete this block to update the Q-Network using the target network
82     1. Compute target using self.Q_target ( target = r + discount * max_b [Q_target(s,b)] )
83     2. Compute Q(s,a) using self.Q
84     3. Compute MSE loss between step 1 and step 2
85     4. Update your network
86     Input: experiences consisting of states,actions,rewards,next_states and discount factor
87     Return: None
88     """
89     states, actions, rewards, next_states, dones = experiences
90     q_eval = self.Q(states).gather(1, actions)
91     q_next = self.Q_target(next_states).detach().max(1)[0].unsqueeze(1)
92     q_target = rewards + discount * q_next * (1-dones)
93     loss = F.mse_loss(q_eval, q_target)
94     self.optimizer.zero_grad()
95     loss.backward()
96     self.optimizer.step()
97     self.target_update(self.Q, self.Q_target, self.tau)
98
99 def target_update(self, Q, Q_target, tau):
100     """
101     TODO: Update the target network parameters (param_target) using current Q parameters (param_Q)
102     Perform the update using tau, this ensures that we do not change the target network drastically
103     1. param_target = tau * param_Q + (1 - tau) * param_target
104     Input: Q,Q_target,tau
105     Return: None
106     """
107     ##### TYPE YOUR CODE HERE #####
108     for param_target, param_local in zip(Q_target.parameters(), Q.parameters()):
109         param_target.data.copy_(tau*param_local.data + (1.0-tau)*param_target.data)
110
111
112     #param_target = tau * param_Q + (1 - tau) * param_target

```

### Training on a fixed epsilon

```

1 def fixed_epsilon_trainer():
2     seed = 0
3     n_episodes = 1000
4     batch_size = 64
5     discount = 0.99
6     lr = 5e-3                # learning rate
7     tau = 0.001              # soft update of target network
8     max_size = int(1e5)
9     update_freq = 4
10    gpu_index = 0
11    max_eps_len = 1000
12    #exploration strategy
13
14    epsilon = 1
15
16    # making the environment
17    env = gym.make('CartPole-v1')
18
19    #setting seeds
20    torch.manual_seed(seed)
21    np.random.seed(seed)
22    random.seed(seed)
23
24    state_dim = env.observation_space.shape[0]
25    action_dim = env.action_space.n
26
27    kwargs = {
28        "state_dim":state_dim,

```

```

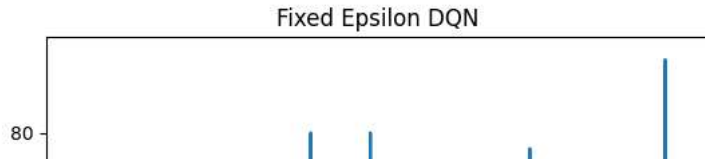
29     "action_dim":action_dim,
30     "discount":discount,
31     "tau":tau,
32     "lr":lr,
33     "update_freq":update_freq,
34
35     "max_size":max_size,
36     "batch_size":batch_size,
37     "gpu_index":gpu_index
38 }
39 fixed_learner = DQNAgent(**kwargs) #Creating the DQN learning agent
40
41 moving_window = deque(maxlen=100)
42
43 reward_store = []
44
45
46 for e in range(n_episodes):
47     state, _ = env.reset(seed=seed)
48     curr_reward = 0
49     for t in range(max_eps_len):
50         action = fixed_learner.select_action(state, epsilon) #To be implemented
51         n_state,reward,terminated,truncated,_ = env.step(action)
52         done = terminated or truncated
53         fixed_learner.step(state, action, reward, n_state, done) #To be implemented
54         state = n_state
55         curr_reward += reward
56         if done:
57             break
58     moving_window.append(curr_reward)
59     reward_store.append(curr_reward)
60
61     """
62     TODO: Write code for decaying the exploration rate using args.epsilon_decay
63     and args.epsilon_end. Note that epsilon has been initialized to args.epsilon_start
64     1. You are encouraged to try new methods
65     """
66     ##### TYPE YOUR CODE HERE #####
67     #####
68
69     if e % 100 == 0:
70         print('Episode Number {} Average Episodic Reward (over 100 episodes): {:.2f}'.format(e, np.mean(moving_window)))
71
72     """
73     TODO: Write code for
74     1. Logging and plotting
75     2. Rendering the trained agent
76     """
77     ##### TYPE YOUR CODE HERE #####
78     #####
79     fig = plt.figure()
80     ax = fig.add_subplot(111)
81     plt.plot(np.arange(len(reward_store)), reward_store)
82     plt.ylabel('Score')
83     plt.xlabel('Episode #')
84     plt.title('Fixed Epsilon DQN')
85     plt.show()
86 fixed_epsilon_trainer()

```

```

Episode Number 0 Average Episodic Reward (over 100 episodes): 18.00
Episode Number 100 Average Episodic Reward (over 100 episodes): 22.46
Episode Number 200 Average Episodic Reward (over 100 episodes): 21.45
Episode Number 300 Average Episodic Reward (over 100 episodes): 21.27
Episode Number 400 Average Episodic Reward (over 100 episodes): 21.25
Episode Number 500 Average Episodic Reward (over 100 episodes): 23.85
Episode Number 600 Average Episodic Reward (over 100 episodes): 20.29
Episode Number 700 Average Episodic Reward (over 100 episodes): 22.21
Episode Number 800 Average Episodic Reward (over 100 episodes): 19.46
Episode Number 900 Average Episodic Reward (over 100 episodes): 23.63

```



```

1 def decaying_epsilon_trainer():
2     seed = 0
3     n_episodes = 1000
4     batch_size = 64
5     discount = 0.99
6     lr = 5e-2                # learning rate
7     tau = 0.001              # soft update of target network
8     max_size = int(1e5)
9     update_freq = 4
10    gpu_index = 0
11    max_eps_len = 1000
12    #exploration strategy
13    # making the environment
14    env = gym.make('CartPole-v1')
15
16    #setting seeds
17    torch.manual_seed(seed)
18    np.random.seed(seed)
19    random.seed(seed)
20    state_dim = env.observation_space.shape[0]
21    action_dim = env.action_space.n
22    kwargs = {
23        "state_dim":state_dim,
24        "action_dim":action_dim,
25        "discount":discount,
26        "tau":tau,
27        "lr":lr,
28        "update_freq":update_freq,
29        "max_size":max_size,
30        "batch_size":batch_size,
31        "gpu_index":gpu_index
32    }
33    decaying_learner = DQNAgent(**kwargs) #Creating the DQN learning agent
34    moving_window = deque(maxlen=100)
35
36    reward_store = []
37
38    epsilon_start = 1          # start value of epsilon
39    epsilon_end = 0.01         # end value of epsilon
40    epsilon_decay = 0.995      # decay value of epsilon
41
42    epsilon_by_step = lambda step: float(epsilon_end+(epsilon_start - epsilon_end)*np.exp(-1. * step / epsilon_decay))
43
44
45    for e in range(n_episodes):
46        state, _ = env.reset(seed=seed)
47        curr_reward = 0
48        for t in range(max_eps_len):
49            action = decaying_learner.select_action(state, epsilon_by_step(t)) #To be implemented
50            n_state,reward,terminated,truncated,_ = env.step(action)
51            done = terminated or truncated
52            decaying_learner.step(state, action, reward, n_state, done) #To be implemented
53            state = n_state
54            curr_reward += reward
55            if np.mean(moving_window) > 200.00:
56                torch.save(agent.Q.state_dict(), 'checkpoint.pth')
57
58            if done:
59                break
60        reward_store.append(curr_reward)

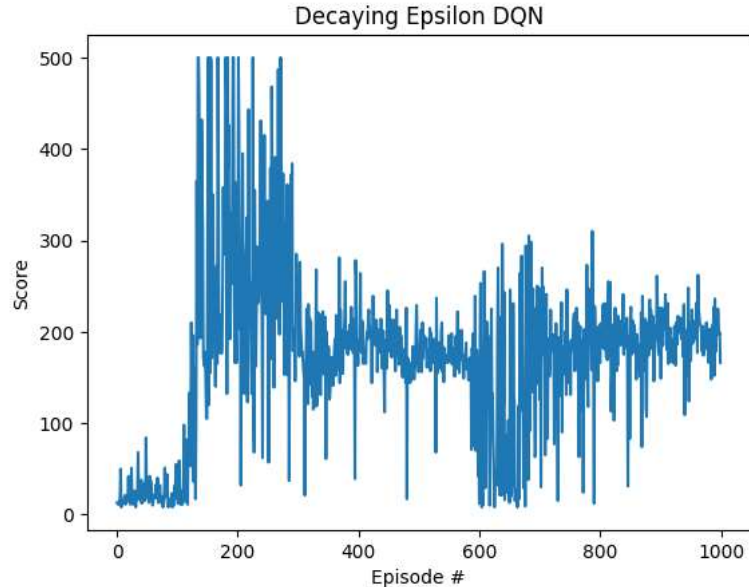
```

```

61 moving_window.append(curr_reward)
62
63 """
64     TODO: Write code for decaying the exploration rate using args.epsilon_decay
65     and args.epsilon_end. Note that epsilon has been initialized to args.epsilon_start
66     1. You are encouraged to try new methods
67 """
68 ##### TYPE YOUR CODE HERE #####
69 #####
70
71
72 if e % 100 == 0:
73     print('Episode Number {} Average Episodic Reward (over 100 episodes): {:.2f}'.format(e, np.mean(moving_window)))
74
75 """
76     TODO: Write code for
77     1. Logging and plotting
78     2. Rendering the trained agent
79 """
80 ##### TYPE YOUR CODE HERE #####
81 #####
82 fig2 = plt.figure()
83 ax = fig2.add_subplot(111)
84 plt.plot(np.arange(len(reward_store)), reward_store)
85 plt.ylabel('Score')
86 plt.xlabel('Episode #')
87 plt.title('Decaying Epsilon DQN')
88 plt.show()
89 decaying_epsilon_trainer()

```

Episode Number 0 Average Episodic Reward (over 100 episodes): 13.00  
 Episode Number 100 Average Episodic Reward (over 100 episodes): 22.18  
 Episode Number 200 Average Episodic Reward (over 100 episodes): 218.99  
 Episode Number 300 Average Episodic Reward (over 100 episodes): 257.91  
 Episode Number 400 Average Episodic Reward (over 100 episodes): 177.86  
 Episode Number 500 Average Episodic Reward (over 100 episodes): 182.68  
 Episode Number 600 Average Episodic Reward (over 100 episodes): 167.80  
 Episode Number 700 Average Episodic Reward (over 100 episodes): 121.35  
 Episode Number 800 Average Episodic Reward (over 100 episodes): 174.20  
 Episode Number 900 Average Episodic Reward (over 100 episodes): 188.14



```

1 # For visualization
2 from gym.wrappers.monitoring import video_recorder
3 from IPython.display import HTML
4 from IPython import display
5 import glob

```

```

1 def show_video(env_name):
2     mp4list = glob.glob('video/*.mp4')
3     if len(mp4list) > 0:
4         mp4 = 'video/{}.mp4'.format(env_name)
5         video = io.open(mp4, 'r+b').read()
6         encoded = base64.b64encode(video)

```

```

7         display.display(HTML(data='''<video alt="test" autoplay
8             loop controls style="height: 400px;">
9                 <source src="data:video/mp4;base64,{0}" type="video/mp4" />
10                </video>'''.format(encoded.decode('ascii'))))
11     else:
12         print("Could not find video")
13
14 def show_video_of_model(agent, env_name):
15     env = gym.make(env_name, render_mode="rgb_array")
16     fourcc = cv2.VideoWriter_fourcc(*'mp4v')
17     video = cv2.VideoWriter('cart_pole.mp4', fourcc, 30, (600, 400))
18     agent.Q.load_state_dict(torch.load('checkpoint.pth'))
19     state, _ = env.reset()
20     print(state.shape)
21     done = False
22     while not done:
23         frame = env.render()
24         video.write(frame)
25
26         action = agent.select_action(state)
27
28         n_state, reward, terminated, truncated, _ = env.step(action)
29         done = terminated or truncated
30         agent.step(state, action, reward, n_state, done) #To be implemented
31         state = n_state
32     env.close()
33     video.release()

```

```

1 import cv2
2 agent = DQNAgent(state_dim=4, action_dim=2)
3 show_video_of_model(agent, 'CartPole')
4

```

```

/usr/local/lib/python3.9/dist-packages/gymnasium/envs/registration.py:531: UserWarning: WARN: Using the latest versioned environment `C
  logger.warn(
(4,)

```

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 11:09 PM

