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.

```
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 16.1 MB/s eta 0:00:00
Requirement already satisfied: numpy>=1.21.0 in /usr/local/lib/python3.9/dist-packages (from gymnasium[box2d]) (1.22.4)
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: 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 farama-notifications>=0.0.1
  Downloading Farama_Notifications-0.0.4-py3-none-any.whl (2.5 kB)
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 58.8 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 14.0 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
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 47.9 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 possib
  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: pyparsing>=2.3.1 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (3.0.9)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (23.0)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (4.39.3)
Requirement already satisfied: importlib-resources>=3.2.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (5.12.0)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (1.0.7)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (0.11.0)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (8.4.0)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.9/dist-packages (from matplotlib) (2.8.2)
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) (3
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.9/dist-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)

```
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 gymnasium.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 dimenssion
 9              action_dim (int): action dimenssion
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)
```

```python
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, Exp_rep, targ_net, 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 netowork
51        """
52        self.t_train += 1
53        self.memory.add(state, action, reward, next_state, done)
54
55
56        # Experience Reply and Target Network Enabled
57        if Exp_rep and targ_net:
58
59            # Conducting Experience Replay
60            if self.memory.size > self.batch_size:
61                experiences = self.memory.sample()
62                self.learn(experiences, self.discount) #To be implemented
63
64                # Updating Target Network
65                if(self.t_train % self.update_freq) == 0:
66                    self.target_update(self.Q, self.Q_target, self.tau)
67
68        # ONLY conducting Experience Replay
69        elif Exp_rep and not targ_net:
70
71            if self.memory.size > self.batch_size:
72                self.memory.add(state, action, reward, next_state, done)
73                experiences = self.memory.sample()
74                self.learn_exp_rep(experiences, self.discount)
75            # Target Network should not be updated
76
77        # ONLY updating Target Network
```

```
 78            else:
 79
 80                if self.memory.size > self.batch_size:
 81                    self.learn_targ_net(state, action, reward, next_state, done, self.discount)
 82                    if(self.t_train % self.update_freq) == 0:
 83                        self.target_update(self.Q, self.Q_target, self.tau)
 84
 85 #To be implemented
 86
 87      def select_action(self, state, epsilon = 0.):
 88          if np.random.random() > epsilon:
 89
 90              state = torch.tensor(state, dtype=torch.float32, device=self.device) # converting our state to pytorch tensor
 91              #self.Q(state) # writing updated tensor to device
 92              self.Q.eval()
 93              with torch.no_grad():
 94                  actions=self.Q.forward(state)
 95              self.Q.train()
 96              action = torch.argmax(actions).item() # item() changes from tensor to integer
 97          else:
 98              action=np.random.choice(self.action_dim)
 99
100          return action
101
102      def learn(self, experiences, discount):
103          """
104          TODO: Complete this block to update the Q-Network using the target network
105          1. Compute target using  self.Q_target ( tar get = r + discount * max_b [Q_target(s,b)] )
106          2. Compute Q(s,a) using self.Q
107          3. Compute MSE loss between step 1 and step 2
108          4. Update your network
109          Input: experiences consisting of states,actions,rewards,next_states and discount factor
110          Return: None
111          """
112          states, actions, rewards, next_states, dones = experiences
113
114
115          q_eval = self.Q(states).gather(1, actions)
116          q_next = self.Q_target(next_states).detach().max(1)[0].unsqueeze(1)
117
118          # Calculating q_target
119          q_target = rewards + discount * q_next * (1-dones)
120
121          # Calculating loss and backpropogating
122          loss = F.mse_loss(q_eval, q_target)
123          self.optimizer.zero_grad()
124          loss.backward()
125          self.optimizer.step()
126          self.target_update(self.Q, self.Q_target, self.tau)
127
128
129      def learn_exp_rep(self, experiences, discount):
130          """
131          TODO: Complete this block to update the Q-Network using the target network
132          1. Compute target using  self.Q_target ( tar get = r + discount * max_b [Q_target(s,b)] )
133          2. Compute Q(s,a) using self.Q
134          3. Compute MSE loss between step 1 and step 2
135          4. Update your network
136          Input: experiences consisting of states,actions,rewards,next_states and discount factor
137          Return: None
138          """
139          states, actions, rewards, next_states, dones = experiences
140
141          q_eval = self.Q(states).gather(1, actions)
142          q_next = self.Q_target(next_states).detach().max(1)[0].unsqueeze(1)
143          q_target = rewards + discount * q_next * (1-dones)
144          loss = F.mse_loss(q_eval, q_target)
145          self.optimizer.zero_grad()
146          loss.backward()
147          self.optimizer.step()
148          #self.target_update(self.Q, self.Q_target, self.tau)
149
150      def learn_targ_net(self, states, actions, rewards, next_states, dones, discount):
151          """
152          TODO: Complete this block to update the Q-Network using the target network
153          1. Compute target using  self.Q_target ( tar get = r + discount * max_b [Q_target(s,b)] )
154          2. Compute Q(s,a) using self.Q
```

```
155          3. Compute MSE loss between step 1 and step 2
156          4. Update your network
157          Input: experiences consisting of states,actions,rewards,next_states and discount factor
158          Return: None
159          """
160          # Storing states
161
162          states, actions, rewards, next_states, dones = [states], [actions], [rewards], [next_states], [dones]
163          states = torch.tensor(states).float().to(self.device)
164          actions = torch.tensor(actions).long().unsqueeze(1).to(self.device)  # Fix here
165          rewards = torch.tensor(rewards).float().unsqueeze(1).to(self.device)  # Fix here
166          next_states = torch.tensor(next_states).float().to(self.device)
167          dones = torch.tensor(dones).unsqueeze(1).float().to(self.device)
168
169          # Calculating q_target
170          q_eval = self.Q(states).gather(1, actions)
171          q_next = self.Q_target(next_states).detach().max(1)[0].unsqueeze(1)
172
173          # Calculating loss and backpropogating
174          q_target = rewards + discount * q_next * (1-dones)
175          loss = F.mse_loss(q_eval, q_target)
176          self.optimizer.zero_grad()
177          loss.backward()
178          self.optimizer.step()
179          self.target_update(self.Q, self.Q_target, self.tau)
180
181      def target_update(self, Q, Q_target, tau):
182          """
183          TODO: Update the target network parameters (param_target) using current Q parameters (param_Q)
184          Perform the update using tau, this ensures that we do not change the target network drastically
185          1. param_target = tau * param_Q + (1 - tau) * param_target
186          Input: Q,Q_target,tau
187          Return: None
188          """
189          ###### TYPE YOUR CODE HERE ######
190          for param_target, param_local in zip(Q_target.parameters(), Q.parameters()):
191              param_target.data.copy_(tau*param_local.data + (1.0-tau)*param_target.data)
192
193
194          #param_target = tau * param_Q + (1 - tau) * param_target
```

Now I'm comparing only target network and Experience Replay. \

This is the code for ONLY Experience Replay

```
1 def exp_rep_only():
2      seed = 0
3      n_episodes = 1500
4      batch_size = 64
5      discount = 0.99
6      lr = 5e-4                         # 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("LunarLander-v2")
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,
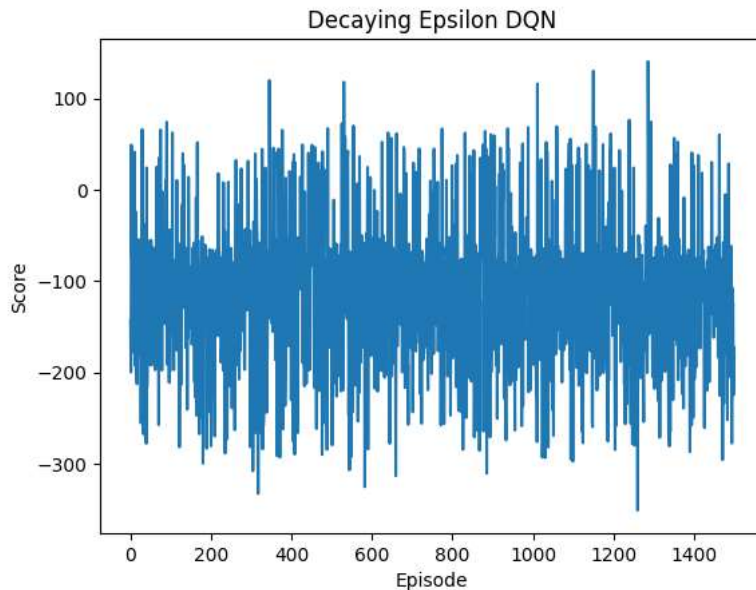```

```
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          exp_rep_only_learner = DQNAgent(**kwargs) #Creating the DQN learning agent
40          moving_window = deque(maxlen=100)
41          reward_store = []
42          print("Experience Replay Only")
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 = exp_rep_only_learner.select_action(state, epsilon) #To be implemented
50                  n_state,reward,terminated,truncated,_ = env.step(action)
51                  done = terminated or truncated
52                  exp_rep_only_learner.step(True, False, state, action, reward, n_state, done) #To be implemented
53                  state = n_state
54                  curr_reward += reward
55                  if done:
56                      break
57              moving_window.append(curr_reward)
58              reward_store.append(curr_reward)
59
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          plt.plot(np.arange(len(reward_store)), reward_store)
69          plt.ylabel('Score')
70          plt.xlabel('Episode')
71          plt.title('Decaying Epsilon DQN')
72          plt.show()
73
74          if e % 100 == 0:
75              print('Episode Number {} Average Episodic Reward (over 100 episodes): {:.2f}'.format(e, np.mean(moving_window)))
76
77
78              ###### TYPE YOUR CODE HERE ######
79              ################################
80          averages = []
81          window_size=50
82
83          for i in range(len(reward_store) - window_size + 1):
84              window = reward_store[i:i+window_size]
85              average = sum(window) / window_size
86              averages.append(average)
87
88          plt.plot(averages, color='g')
89          plt.ylabel('Score')
90          plt.xlabel('Episodes')
91          plt.title('Only Experience Replay DQN')
92          plt.show()
93
94 exp_rep_only()
```

Experience Replay Only

### Decaying Epsilon DQN



### Only Experience Replay DQN



This is the code for only Target Network

```
1   def only_tar_net_trainer():
2       seed = 0
3       n_episodes = 500
4       batch_size = 64
5       discount = 0.99
6       lr = 5e-4                          # 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
13      #exploration strategy
14
15      # making the environment
16      env = gym.make("LunarLander-v2")
17      #setting seeds
18      torch.manual_seed(seed)
19      np.random.seed(seed)
20      random.seed(seed)
21      state_dim = env.observation_space.shape[0]
22      action_dim = env.action_space.n
23      kwargs = {
24          "state_dim":state_dim,
25          "action_dim":action_dim,
26          "discount":discount,
27          "tau":tau,
28          "lr":lr,
29          "update_freq":update_freq,
30          "max_size":max_size,
31          "batch_size":batch_size,
32          "gpu_index":gpu_index
33      }
34      disabled_decaying_learner = DQNAgent(**kwargs) #Creating the DQN learning agent
35      moving_window = deque(maxlen=100)
36
37      reward_store = []
38
```
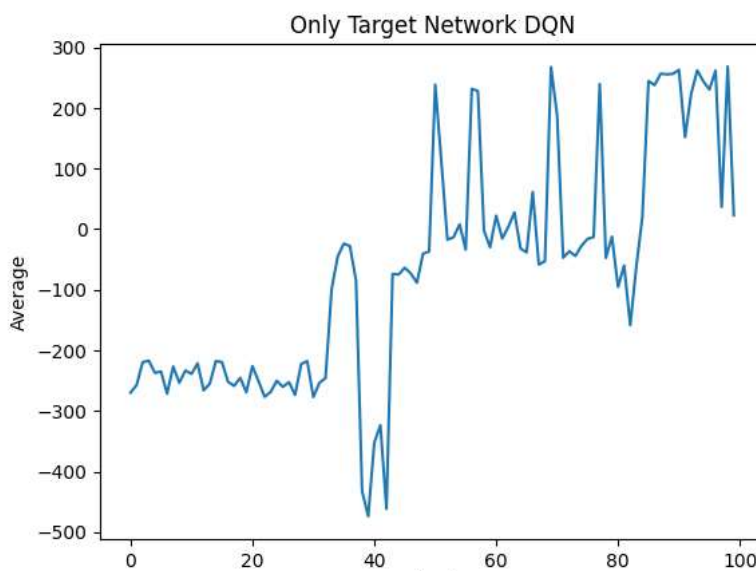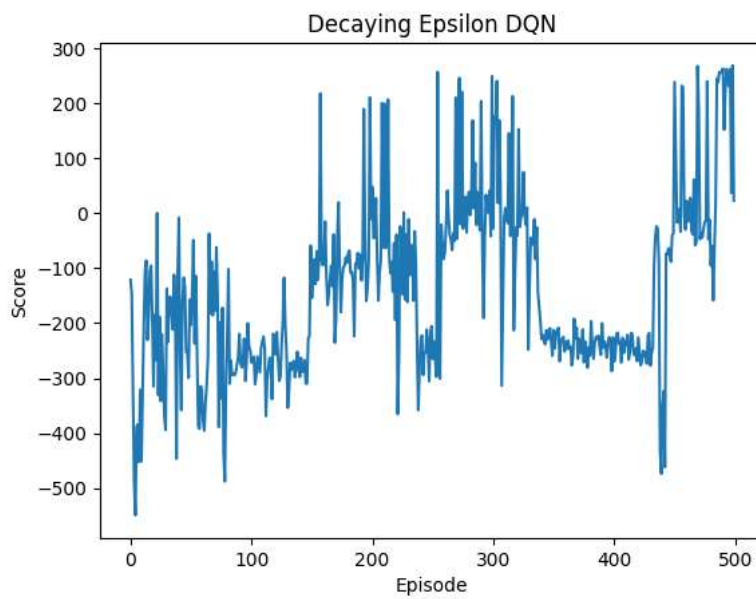
```
39       print("Target Network DQN")
40
41       epsilon_start = 1                  # start value of epsilon
42       epsilon_end = 0.01                 # end value of epsilon
43       epsilon_decay = 0.995              # decay value of epsilon
44
45       epsilon_by_step = lambda step: float(epsilon_end+(epsilon_start - epsilon_end)*np.exp(-1. * step / epsilon_decay))
46
47
48       for e in range(n_episodes):
49         state, _ = env.reset(seed=seed)
50         curr_reward = 0
51         for t in range(max_eps_len):
52           action = disabled_decaying_learner.select_action(state, epsilon_by_step(t)) #To be implemented
53           n_state,reward,terminated,truncated,_ = env.step(action)
54           done = terminated or truncated
55           disabled_decaying_learner.step(False, True, state, action, reward, n_state, done) #To be implemented
56           state = n_state
57           curr_reward += reward
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         if e % 100 == 0:
72           print('Episode Number {} Average Episodic Reward (over 100 episodes): {:.2f}'.format(e, np.mean(moving_window)))
73
74         """
75           TODO: Write code for
76           1. Logging and plotting
77           2. Rendering the trained agent
78         """
79         ###### TYPE YOUR CODE HERE ######
80         ################################
81       plt.plot(np.arange(len(reward_store)), reward_store)
82       plt.ylabel('Score')
83       plt.xlabel('Episode')
84       plt.title('Decaying Epsilon DQN')
85       plt.show()
86
87       plt.plot(np.arange(len(moving_window)), moving_window)
88       plt.ylabel('Average')
89       plt.xlabel('Episodes')
90       plt.title('Only Target Network DQN')
91       plt.show()
92
93   only_tar_net_trainer()
```

```
Target Network DQN
Episode Number 0 Average Episodic Reward (over 100 episodes): -121.30
Episode Number 100 Average Episodic Reward (over 100 episodes): -246.68
Episode Number 200 Average Episodic Reward (over 100 episodes): -170.87
Episode Number 300 Average Episodic Reward (over 100 episodes): -64.01
Episode Number 400 Average Episodic Reward (over 100 episodes): -152.57
```

Decaying Epsilon DQN

Only Target Network DQN

Colab paid products  -  Cancel contracts here

✓  6m 12s  completed at 11:03 PM