

COURSEWORK 2

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

COMP97143: Reinforcement Learning

Author:

Jaime Sabal Bermúdez (CID: 01520988)

Date: November 29, 2021

Question 1: Implementing a Functional DQN

1.1-1.2: Implementation of DQN

Replay Buffer).

The `ReplayBuffer` class was implemented with a double-ended queue (i.e deque) as its main component (line 26 in source code) which allowed for an easy way to keep track a number of previously sampled states (given by the variable `buffer_size`). After the agent performs a step in the environment, we make use of the `push` method (line 28) to save this transition. This transition is a tuple of the form (s_t, a_t, s_{t+1}, r_t) , where s_t is the current state, a_t is the action taken, s_{t+1} is the next state and r_t is the reward obtained from transitioning from s_t to s_{t+1} . Once enough states have been stored in the memory (more than the batch size of the Deep Q Network (DQN)), the weights of the network are optimised by randomly sampling a batch of previously seen states. Doing this allows for less correlation between training samples and for these to follow a similar probability distribution (i.i.d) like is expected in supervised learning.

Target Network).

The target network is initialised in lines 118-119 of the source code by copying the states dictionary of the policy network. This network is used in calculating the Temporal Difference (TD) error, which stabilises the learning process by allowing the policy network to consider a variety of actions before updating its weights according to an error that may not be representative of the optimum actions to take for each state.

Frame Stacking and Skipping).

We also implement the feature of being able to add multiple frames k as input to the DQN such that it contains the states $(s_t, s_{t-1}, \dots, s_{t-k})$ and thus has shape $(1, 4k)$. We create the class `FrameStacking` (line 57), which inherits from `gym.Wrapper`, with a `reset_buffer` method that empties the memory deque and a `k_states` property that converts the deque to the expected form of the DQN input. We also implement frame skipping (lines 233-243), mostly to accelerate training times when performing replications of our experiments, since the dynamics of the system are set such that one step corresponds to 0.02 seconds and that is close to the average 0.025 second reaction time of a human.

1.3: Justification of Network Architecture

Figure 1 shows the chosen architecture of our DQN. Four frames ($k = 4$), each containing four observables $(x, \dot{x}, \theta, \dot{\theta})$, were used as input to the network, 150 neurons are present in each hidden layer, and the ReLu activation function operates on all layers but the output layer. Moreover, the Adam optimiser with a learning rate of $\alpha = 0.0001$ was used to update the weights in training with minibatches of size 128, which we found yielded better results than using RMSProp and larger values for α after some experimentation.

In terms of hyperparameters specific to a DQN, we use a discount rate $\gamma = 0.9999$, update the target network every 40 episodes, use a replay buffer size of 100,000, and a rewards-based ϵ -greedy behaviour policy with a starting $\epsilon_i = 1.0$ and a final $\epsilon_f = 0.0005$. Specifically, in a rewards-based ϵ -decay ϵ is reduced by $\delta = (\epsilon_i - \epsilon_f)/r_T$ every time the reward obtained in an episode is larger than a threshold

value that is initialised to be 0 and increased by 1 every time this condition is met.

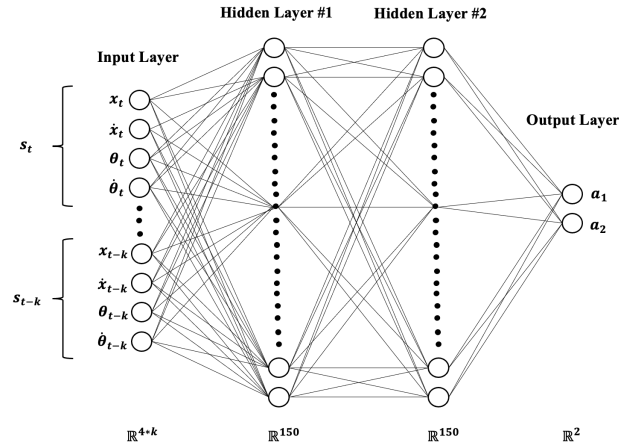


Figure 1: Graphical representation of the architecture used in the Deep Q Network. Each connection to the center of the hidden layers represents the neurons' connection to each of the remaining neurons (i.e 146) in the subsequent layer.

In the mentioned equation, r_T is a reward target value that we set to be 100. By doing this, the decay is much more controlled and only decreases if our agent obtains high rewards, which in itself is an indication that our agent should explore less. Moreover, the discount rate γ is set very high due to the fact that we are interested in obtaining high rewards in the very far future. Finally, updating the target network every 40 episodes (f_T) and using a replay buffer size of 100,000 allowed for a more stable network that consistently converged to a reward of ≈ 300 in 700 episodes without suffering from catastrophic forgetting.

1.4: Learning Curve of Agent

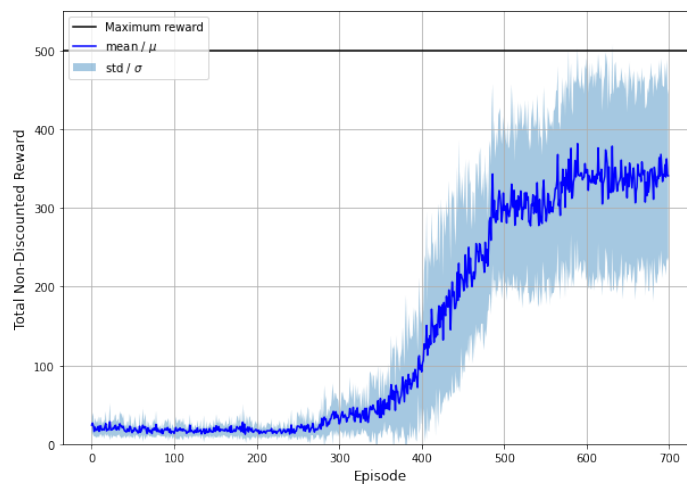


Figure 2: Average learning curve and its standard deviation of a DQN agent solving OpenAI's cartpole environment across 15 replications; $\text{batch_size} = 128$, $\alpha = 0.0001$, $f_T = 40$, $k = 4$, $\text{skipped_frames} = 4$, $\text{buffer_size} = 100,000$, $\epsilon_i = 1.0$, $\epsilon_f = 0.0005$, $r_T = 100$.

Figure 2 shows the average learning curve of the DQN agent across 15 replications for 700 training episodes. The range chosen for the plot is $[0, 550]$ such that it clearly shows the maximum attainable reward (500) for any given episode. This number of replications was chosen based on an informal analysis of the tradeoff between the marginal benefit of an extra replication and the time taken to compute it. Specifically, after 15 replications the smoothness of the average reward and standard deviation curves didn't seem to improve much. Moreover, the average total return of the final 100 episodes during training was found to be ≈ 339 , resulting in our agent achieving 90% of this expected performance (305) at episode 485.

Question 2: Hyperparameters of the DQN

2.1: Effect of ϵ on the Learning Curve of the Agent

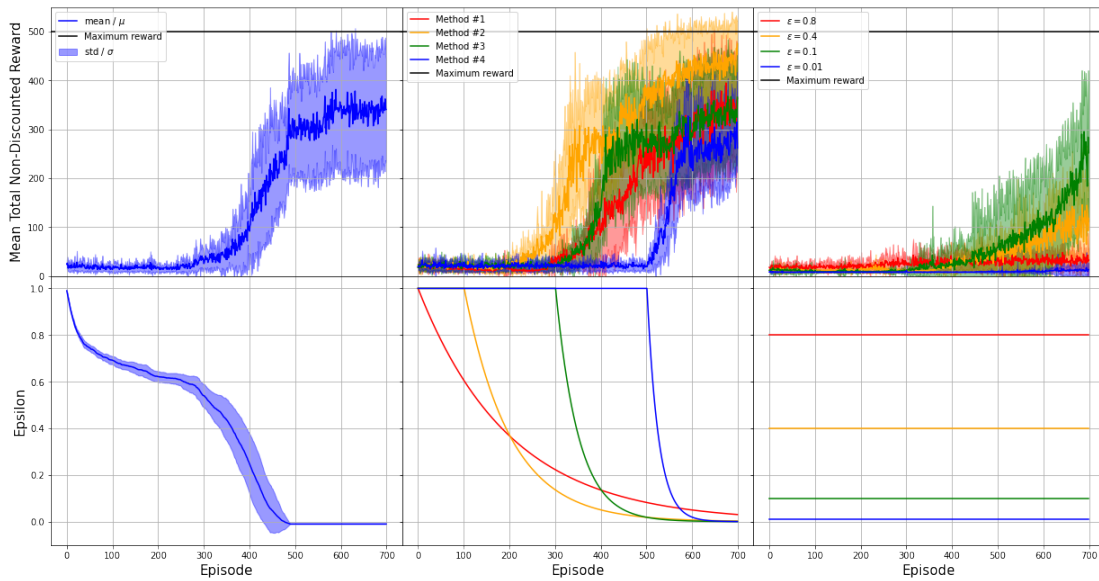


Figure 3: Learning curves of DQN agent for different ϵ schedules with the same starting and ending $\epsilon = 1.0$ and $\epsilon = 0.0005$, respectively; **left**). reward-based decay with a rewards target of 100 (also shows σ for the variation in the decay for different runs); **center**). Exponential ϵ decay with the decay beginning at different episodes in training for different methods; #1 = 0, #2 = 100, #3 = 300, #4 = 500; **right**). Constant ϵ throughout learning.

Figure 3 shows the learning curves for the DQN agent using different ϵ schedules while keeping the other hyper-parameters constant and equal to those detailed in Question 1. As can be seen, the best schedules seem to be the rewards-based and exponential decays (left and center columns). The exponential decay schedules were done by starting the attenuation at different episodes following the equation:

$$\epsilon_t = \epsilon_f + (\epsilon_i - \epsilon_f) \exp\left[\frac{-t}{\tau}\right] \quad (1)$$

where t is the episode and τ is a decay factor. The different curves shown use $\tau = 200$ (#1), $\tau = 100$ (#2), $\tau = 50$ (#3), $\tau = 25$ (#4) and begin the attenuation in episodes 0, 100, 300, and 500, respectively. This plot shows how curve #2 is best for this

schedule, allowing for some exploration at the start but giving the agent enough time to learn the environment very well in just 700 episodes. For a constant epsilon throughout training (right-most column), the agent is unable to learn anything for very large or small ϵ since choosing a random action way too often or virtually ever inhabilitates learning (agent must do some exploring but not too much). However, for the intermediate $\epsilon = 0.1$ the agent seems to start learning at a very late episode, which makes sense for a policy that takes a long time to sample all possible scenarios in the environment. Therefore, an ϵ that varies during training is preferred but not completely necessary, since using an appropriate constant ϵ still allows for some degree of learning.

2.2: Effect of Replay Buffer on Variability of Reward During Learning

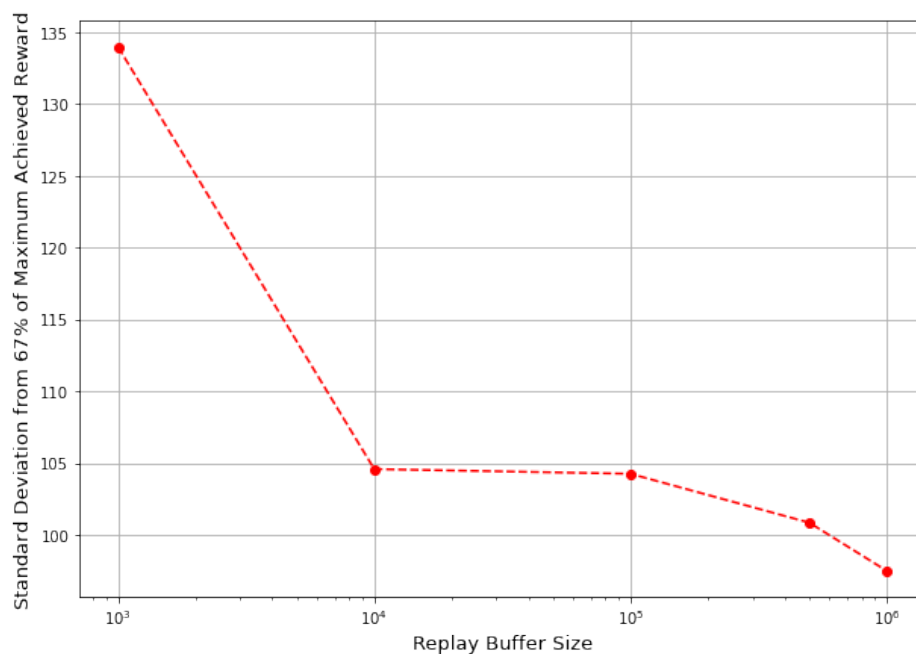


Figure 4: Replay buffer size (in a log scale) against variability in learning curve from 67% of the maximum achieved total reward. The remaining hyper-parameters are the same as those used to compute Figure 2

Exploring the effect of the replay buffer size on the total reward can give us insight on the importance of experience replay in a Deep Q Network. In Figure 4, it can be seen how there is a clear trend that as the buffer size is increased, the variability in the reward during learning from 67% of the maximum achieved reward decreases. One reason for this is that as the agent keeps track of more experiences from the past, the sampled data at each step of the episode is less correlated and the agent is thus able to converge better to a policy that yields consistently good rewards since the data is more independent and identically distributed (i.i.d) like that assumed in supervised learning tasks. In other words, by having a larger amount of previous episodes available to sample from, the agent is less likely to depend on any given run and thus stability is increased in learning.

2.3: Effect of k on the Overall Learning and Final Performance

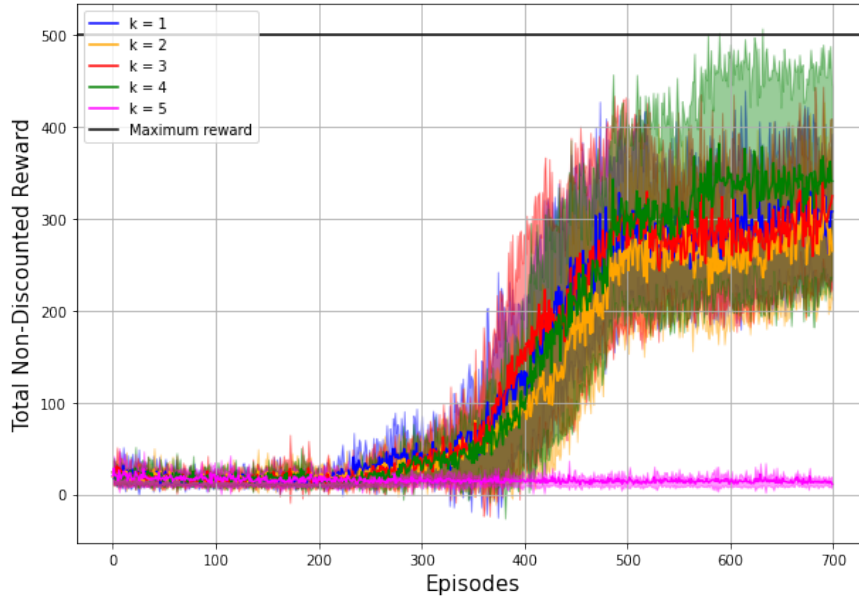


Figure 5: Learning curves for different numbers of frames $k \in [1, 5]$ used as input to the policy network in determining the optimum action to take at any given state. The remaining hyper-parameters used are the same as those shown on Figure 2.

The upper bound for k should be equal to the number of actions after which our agent would not be able to recover if it were to perform the same action repeatedly from the initial state with $\theta = 0$. After performing 100,000 runs of the environment with this policy, I found that the average number of steps needed to reach a terminal state when performing a constant action (either always 0 or always 1) is 9.35. This implies that after performing 5 steps in the same direction our agent will, on average, not be able to perform enough steps in the other direction to recover. Hence, using a $k > 5$ will impede the agents learning process because the network's input will always contain at least one frame that contains conflicting information on the optimum action to take at the current state. Figure 5 shows this phenomenon clearly, along with four other learning curves for $1 \leq k < 5$. For the latter, no major improvement in either the maximum achieved rewards or the variability of the learning curves can be seen. This is presumably due to the environment being deterministic and Markovian, implying that only the current state is necessary for determining the action that will result in the largest total non-discounted rewards. This means that the extra states being inputted to the network in the case where $k > 1$ are redundant and don't bear importance in the agent's learning.

Question 3: Ablation/Augmentation Experiments

3.1: DDQN Implementation

We implement the Double Deep Q Network agent by modifying the DQN agent slightly. In DQN, we calculate the target used in computing the MSE loss of our network (expected state-action value function) by:

$$Y_t^{DQN} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta'_t) \quad (2)$$

where θ'_t are the weights of the target network and we are thus using the same network to both select and evaluate any given action. Since the value function is prone to be very noisy, taking its maximum leads to overestimation errors that dampen the learning process (maximisation bias). This occurs because the values for different actions are most likely not overestimated by equal amounts and thus for any given episode the weights of the behaviour policy θ_t may be updated to favour actions that are in reality not optimal. The solution to this problem is to select the actions with the online behaviour policy and use these actions to evaluate the target network (Double Q Network). This can be written as:

$$Y_t^{DDQN} = R_{t+1} + \gamma Q\left(S_{t+1}, \operatorname{argmax}_a Q(S_{t+1}, a; \theta_t); \theta'_t\right) \quad (3)$$

The idea behind this method to calculate the TD errors is that because the weights θ'_t are updated on a periodic basis (every 40 episodes in our case), the likeliness of both policies overestimating the same action is greatly reduced since the noise present in both of them are independent of each other. The implementation of this augmentation to the idea of a target network can be seen between lines 181 and 185 of the source code.

3.2: Learning Curves for DQN Modifications Against Original Model

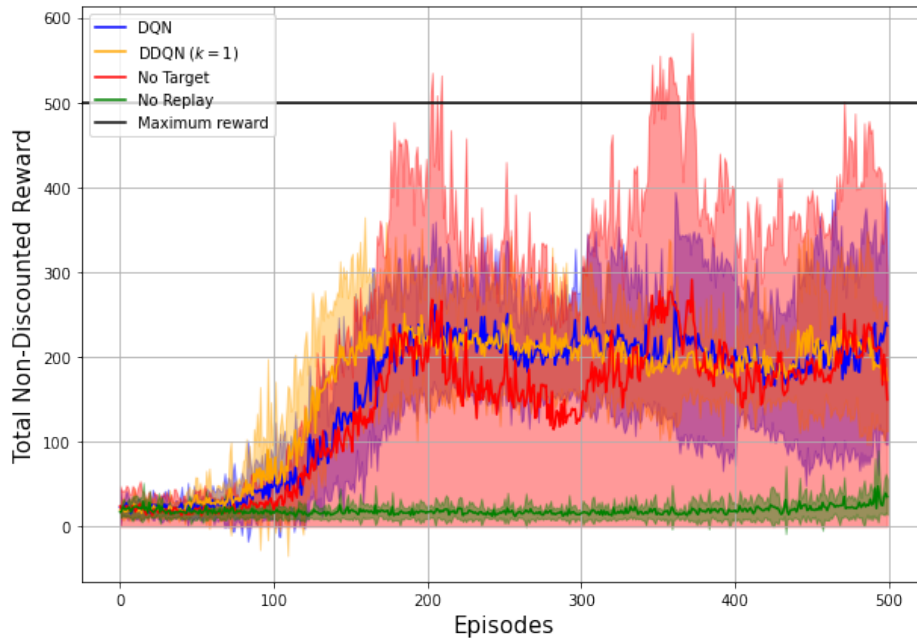


Figure 6: Learning curves for a DQN agent, one without a target network, one without experience replay, and a DDQN agent. `batch_size = 128` (for No Replay `batch_size = 1`), $\alpha = 0.0001$, $f_T = 40$ (except for No Target), $k = 4$ (except for DDQN where $k = 1$), `skipped_frames = 1`, `buffer_size = 100,000`, $\epsilon_i = 1.0$, $\epsilon_f = 0.0005$, $r_T = 100$.

As can be seen in Figure 6, the learning curves for the two ablations is quite different from that of the DQN and DDQN agents. It is important to note that for this section of the report, we eliminated frame skipping due to time constraints because it allowed for convergence in less episodes, albeit this also resulted in a lower final performance of the agents (see Figure 2).

Without having experience replay, the agent forgets previous transitions and is thus unable to learn from its mistakes. On the other hand, removing the target network has the effect of producing a very unstable learning process due to the fact that we are bootstrapping a continuous state space representation. In other words, updating the Q network for one state updates it for all the other states too which can lead to a resonance effect (i.e catastrophic forgetting). This is evidenced by the oscillating nature of the curve (red) and the absurdly high variance seen between replications.

Finally, the DDQN model shows a very similar learning curve to the DQN agent. However, a key difference between the two is the time taken to converge to a solution. The DDQN's learning curve can be seen to reach a similar final performance more quickly. This is due to the agent spending less time following un-optimal policies, since the Q-values of these state-action pairs are much less likely to be larger than the optimal alternative.

1 Appendix

1.1 Source Code

```
1 import gym
2 import copy
3 from gym.wrappers.monitoring.video_recorder import VideoRecorder #
  ↳ records videos of episodes
4
5 import numpy as np
6 import matplotlib.pyplot as plt # Graphical library
7 import torch
8 import torch.optim as optim
9 import torch.nn as nn
10 import torch.nn.functional as F
11 import tqdm
12
13 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
  ↳ # Configuring Pytorch
14 from collections import namedtuple, deque
15 from itertools import count
16 import math
17 import random
18
19 Transition = namedtuple('Transition',
20                        ('state', 'action', 'next_state', 'reward'))
21
22 #dont change
23 class ReplayBuffer(object):
24
25     def __init__(self, capacity):
26         self.memory = deque([],maxlen=capacity)
27
28     def push(self, *args):
29         """Save a transition"""
30         self.memory.append(Transition(*args))
31
32     def sample(self, batch_size):
33         return random.sample(self.memory, batch_size)
34
35     def __len__(self):
36         return len(self.memory)
37
38 #dont change
39 class DQN(nn.Module):
40
```

```

41     def __init__(self, inputs, outputs, num_hidden, hidden_size):
42         super(DQN, self).__init__()
43         self.input_layer = nn.Linear(inputs, hidden_size)
44         self.hidden_layers = nn.ModuleList([nn.Linear(hidden_size,
45             ↪ hidden_size) for _ in range(num_hidden-1)])
46         self.output_layer = nn.Linear(hidden_size, outputs)
47
48     def forward(self, x):
49         x.to(device)
50
51         x = F.relu(self.input_layer(x))
52         for layer in self.hidden_layers:
53             x = F.relu(layer(x))
54
55         return self.output_layer(x)
56
57 class FrameStacking(gym.Wrapper):
58     """Return only every 4th frame"""
59     def __init__(self, env, k):
60         super(FrameStacking, self).__init__(env)
61         self._obs_buffer = deque([], maxlen=k)
62         self._skip      = k
63
64     @property
65     def _k_states(self):
66         return torch.stack(list(self._obs_
67             ↪ buffer)).reshape(-1).unsqueeze(0)
68
69     def reset_buffer(self):
70         self._obs_buffer = deque([], maxlen=self._skip)
71
72 class CartpoleAgent():
73     def __init__(self, NUM_EPISODES, BATCH_SIZE, GAMMA,
74         ↪ TARGET_UPDATE_FREQ,
75             NUM_FRAMES, SKIPPED_FRAMES, NUM_HIDDEN_LAYERS,
76             ↪ SIZE_HIDDEN_LAYERS, REPLAY_BUFFER, LR,
77             EPS_START, EPS_END, REWARD_TARGET, EPS_DECAY=None,
78             ↪ DDQN=False, eps_decay_strat="reward",
79             ablate_target=False, ablate_replay=False):
80
81         #initialisation attributes
82         self.reward_target = REWARD_TARGET
83         self.reward_threshold = 0

```

```
81     self.num_episodes = NUM_EPISODES
82     self.batch_size = BATCH_SIZE
83     self.gamma = GAMMA
84     self.target_update_freq = TARGET_UPDATE_FREQ
85     self.eps_start = EPS_START
86     self.eps_end = EPS_END
87     self.k = NUM_FRAMES
88     self.skip_frames = SKIPPED_FRAMES
89     self.curr_episode = 0
90     self.losses = []
91     self.train_rewards = []
92     self.eps_decay_strat = eps_decay_strat
93
94     self.ablate_target = ablate_target
95     self.ablate_replay = ablate_replay
96
97     if self.eps_decay_strat[:3] == "exp":
98         if EPS_DECAY == None:
99             raise Exception("Have to define EPS_DECAY if decay
100                             ↪ strategy is exponential")
101         else:
102             self.eps_decay = EPS_DECAY
103
104     self.epsilon = self.eps_start
105     self.epsilon_list = []
106     self.epsilon_delta = (self.epsilon -
107                             ↪ self.eps_end)/self.reward_target
108
109     #Get number of states and actions from gym action space
110     env = gym.make("CartPole-v1")
111     env.reset()
112     self.state_dim = len(env.state)    #x, x_dot, theta, theta_dot
113     self.n_actions = env.action_space.n
114     env.close()
115
116     self.input_dim = int(self.state_dim*self.k) #define input size
117     ↪ of DQN
118
119     #define policy and target networks, as well as optimizer and
120     ↪ replay buffer
121     self.policy_net = DQN(self.input_dim, self.n_actions,
122                             ↪ NUM_HIDDEN_LAYERS, SIZE_HIDDEN_LAYERS).to(device)
123     self.target_net = DQN(self.input_dim, self.n_actions,
124                             ↪ NUM_HIDDEN_LAYERS, SIZE_HIDDEN_LAYERS).to(device)
125     self.target_net.load_state_dict(self.policy_net.state_dict())
```

```

120     self.target_net.eval()
121
122     self.DDQN = DDQN
123     self.optimizer = optim.Adam(self.policy_net.parameters(),
124     ↪ lr=LR)
125     self.memory = ReplayBuffer(REPLAY_BUFFER)
126
127
128
129 def select_action(self, k_states):
130     """
131     sample = random.random()
132     if sample > self.epsilon:
133         with torch.no_grad():
134             # t.max(1) will return largest column value of each
135             ↪ row.
136             # second column on max result is index of where max
137             ↪ element was
138             # found, so we pick action with the larger expected
139             ↪ reward.
140             return self.policy_net(k_states).max(1)[1].view(1, 1)
141     else:
142         return torch.tensor([[random.randrange(self.n_actions)]],
143         ↪ device=device, dtype=torch.long)
144
145
146 def optimize_model(self):
147     if len(self.memory) < self.batch_size:
148         return
149     transitions = self.memory.sample(self.batch_size)
150     # Transpose the batch (see
151     ↪ https://stackoverflow.com/a/19343/3343043 for
152     # detailed explanation). This converts batch-array of
153     ↪ Transitions
154     # to Transition of batch-arrays.
155     batch = Transition(*zip(*transitions))
156
157     # Compute a mask of non-final states and concatenate the
158     ↪ batch elements
159     # (a final state would've been the one after which simulation
160     ↪ ended)
161     non_final_mask = torch.tensor(tuple(map(lambda s:
162     ↪ torch.sum(s[0][-self.k:].absolute().item() > 0,
163     ↪ batch.next_state)), device=device, dtype=torch.bool)

```

```

154
155     # Can safely omit the condition below to check that not all
156     ↳ states in the
157     # sampled batch are terminal whenever the batch size is
158     ↳ reasonable and
159     # there is virtually no chance that all states in the sampled
160     ↳ batch are
161     # terminal
162     if sum(non_final_mask) > 0:
163         non_final_next_states = torch.cat([s for s in
164             ↳ batch.next_state if
165             ↳ torch.sum(s[0][-self.k:]).absolute().item() > 0])
166     else:
167         non_final_next_states =
168             ↳ torch.empty(0, self.state_dim).to(device)
169
170     state_batch = torch.cat(batch.state)
171     action_batch = torch.cat(batch.action)
172     reward_batch = torch.cat(batch.reward)
173
174     # Compute Q(s_t, a) - the model computes Q(s_t), then we
175     ↳ select the
176     # columns of actions taken. These are the actions which
177     ↳ would've been taken
178     # for each batch state according to policy_net
179     state_action_values = self.policy_net(state_batch).gather(1,
180         ↳ action_batch)
181
182     # Compute V(s_{t+1}) for all next states.
183     # This is merged based on the mask, such that we'll have
184     ↳ either the expected
185     # state value or 0 in case the state was final.
186     next_state_values = torch.zeros(self.batch_size,
187         ↳ device=device)
188
189     with torch.no_grad():
190         # Once again can omit the condition if batch size is
191         ↳ large enough
192         if sum(non_final_mask) > 0:
193             if self.DDQN:
194                 #DDQN ---> update next states Q values (using
195                 ↳ target net) using the actions that maximise
196                 ↳ the policy network
197                 actions_non_final = torch.zeros_like(action_
198                     ↳ batch.view(non_final_mask.shape))

```

```

184         actions_non_final = torch.argmax(self.policy_net(
            ↪ net(non_final_next_states),
            ↪ 1).unsqueeze(1)
185         next_state_values[non_final_mask] = self.target_net(
            ↪ net(non_final_next_states).gather(1,
            ↪ actions_non_final).flatten()
186     else:
187         if self.ablate_target:
188             next_state_values[non_final_mask] =
            ↪ self.policy_net(non_final_next_
            ↪ states).max(1)[0].detach()
189         else:
190             next_state_values[non_final_mask] =
            ↪ self.target_net(non_final_next_
            ↪ states).max(1)[0].detach()
191
192     else:
193         next_state_values =
            ↪ torch.zeros_like(next_state_values)
194
195
196     expected_state_action_values = (next_state_values *
            ↪ self.gamma) + reward_batch
197
198     # Compute loss
199     loss_func = nn.MSELoss()
200     loss = loss_func(state_action_values,
            ↪ expected_state_action_values.unsqueeze(1))
201
202
203     # Optimize the model
204     self.optimizer.zero_grad()
205     loss.backward()
206
207     # Limit magnitude of gradient for update step
208     #for param in self.policy_net.parameters():
209     #     param.grad.data.clamp_(-1, 1)
210
211     self.optimizer.step()
212
213 def train(self):
214     """ """
215     env = gym.make("CartPole-v1")
216     steps = 0
217     stack_frames = FrameStacking(env, self.k) #define frame
            ↪ stacking framework

```

```
218     decay_iter = 0
219
220     for i_episode in tqdm.tqdm(range(self.num_episodes)):
221         rewards = 0
222         #if i_episode % 20 == 0:
223             #print("episode ", i_episode, "/", self.num_episodes)
224
225         env.reset() #reset environment
226         state =
227             ↪ torch.tensor(env.state).float().unsqueeze(0).to(device)
228         for _ in range(self.k): # Added
229             stack_frames._obs_buffer.append(state) # Added
230
231         for t in count():
232             k_states = stack_frames._k_states
233             temp_rewards = 0
234             for _ in range(self.skip_frames):
235                 #process frames to pass as input to DQN
236                 action = self.select_action(k_states)
237                 _, reward, done, _ = env.step(action.item()) #take
238                 ↪ step following epsilon-greedy policy
239
240             if not done:
241                 k_states = stack_frames._k_states
242             else:
243                 break
244
245             temp_rewards += reward
246
247         rewards += temp_rewards
248         reward = torch.tensor([temp_rewards], device=device)
249
250         # Observe new state
251         if not done:
252             next_state =
253                 ↪ torch.tensor(env.state).float().unsqueeze(0).to(device)
254         else:
255             next_state = torch.zeros(1, self.state_dim)
256
257         #store new state in frames buffer and process frames
258             ↪ to pass as input to DQN
259         stack_frames._obs_buffer.append(next_state)
260         next_k_states = stack_frames._k_states
261
262         #Store the transition in memory
```

```
259         self.memory.push(k_states, action, next_k_states,
260                             ↪ reward)
261
262         # Perform one step of the optimization (on the policy
263         ↪ network)
264         self.optimize_model()
265
266         if done:
267             break
268
269         # Select and perform an action
270         steps += 1
271
272     stack_frames.reset_buffer() #reset frame stack memory
273
274     # Update the target network, copying all weights and
275     ↪ biases in DQN
276     if i_episode % self.target_update_freq == 0:
277         self.target_net.load_state_dict(self.policy_net.state_
278                                         ↪ dict())
279
280     #epsilon decay strategy
281     if self.eps_decay_strat == "reward":
282         #implement reward-based decay
283         if self.epsilon > self.eps_end and rewards >
284             ↪ self.reward_threshold:
285             self.epsilon -= self.epsilon_delta
286             self.reward_threshold += 1
287
288     elif self.eps_decay_strat == "exp1":
289         #epsilon starts decaying exponentially at start of
290         ↪ training
291         self.epsilon = self.eps_end +
292             ↪ (self.eps_start-self.eps_end)*math.exp(-i_
293             ↪ episode/self.eps_decay)
294
295     elif self.eps_decay_strat == "exp2":
296         #epsilon starts decaying exponentially after the
297         ↪ 100th episode
298         if i_episode > 100:
299             self.epsilon = self.eps_end +
300                 ↪ (self.eps_start-self.eps_
301                 ↪ end)*math.exp(-2*decay_iter/self.eps_decay)
302             decay_iter += 1
303
304     elif self.eps_decay_strat == "exp3":
305         #epsilon starts decaying exponentially after the
306         ↪ 300th episode
```



```
293         if i_episode > 300:
294             self.epsilon = self.eps_end +
                ↪ (self.eps_start-self.eps_
                ↪ end)*math.exp(-4*decay_iter/self.eps_decay)
295             decay_iter += 1
296         elif self.eps_decay_strat == "exp4":
297             #epsilon starts decaying exponentially after the
                ↪ 500th episode
298             if i_episode > 500:
299                 self.epsilon = self.eps_end +
                    ↪ (self.eps_start-self.eps_
                    ↪ end)*math.exp(-8*decay_iter/self.eps_decay)
300                 decay_iter += 1
301         else:
302             pass
303
304         self.epsilon_list.append(self.epsilon)
305         self.train_rewards.append(rewards) #append episode reward
                ↪ to list
306
307         self.curr_episode += 1
308
309     print("Complete")
310     env.close() #close environment
311
312
313     def test(self):
314         """run an episode with trained agent and record video"""
315
316         env = gym.make("CartPole-v1")
317         file_path = 'video.mp4'
318         recorder = VideoRecorder(env, file_path)
319
320         env.reset()
321         done = False
322
323         stack_frames = FrameStacking(env, self.k) #define frame
                ↪ stacking framework
324
325         state =
                ↪ torch.tensor(env.state).float().unsqueeze(0).to(device)
326         for _ in range(self.k): # Added
327             stack_frames._obs_buffer.append(state) # Added
328
329         duration = 0
```

```

330
331     while not done:
332         recorder.capture_frame()
333         # Select and perform an action
334         k_states = stack_frames._k_states
335         action = self.select_action(k_states)
336         _, reward, done, _ = env.step(action.item()) #take step
337         ↳ following epsilon-greedy policy
338         duration += 1
339
340         next_state =
341         ↳ torch.tensor(env.state).float().unsqueeze(0).to(device)
342
343         #store new state in frames buffer and process frames to
344         ↳ pass as input to DQN
345         stack_frames._obs_buffer.append(next_state)
346
347     env.close()
348     recorder.close()
349
350     print("Episode duration: ", duration)
351
352 def plot_rewards(self):
353     """Plot the total non-discounted sum of rewards across the
354     ↳ episodes (i.e duration of each episode in steps)."""
355     epochs = np.linspace(0, len(self.train_rewards),
356     ↳ len(self.train_rewards))
357
358     plt.figure(figsize=(10,7))
359     plt.grid()
360     plt.plot(epochs, self.train_rewards, label="Total
361     ↳ Non-Discounted Rewards")
362     #plt.axhline(y=self.reward_target, color='black',
363     ↳ linestyle='-', label="Reward target")
364     #plt.plot(epochs, rewards_list, "b.", markersize=3)
365     #plt.plot(epochs, np.poly1d(np.polyfit(epochs, rewards_list,
366     ↳ 1))(epochs), "r")
367     plt.xlabel("Episodes")
368     plt.ylabel("Total Non-Discounted Reward")
369     plt.ylim((0,550))
370     plt.axhline(y=500, color='red', linestyle='-', label="Maximum
371     ↳ reward")
372     plt.legend(loc="best")
373     plt.show()

```

```

366
367 def plot_losses(self):
368     """Plot the total non-discounted sum of rewards across the
369     ↪ episodes (i.e duration of each episode in steps)."""
370     epochs = np.linspace(0, len(self.losses), len(self.losses))
371
372     plt.figure(figsize=(10,7))
373     plt.grid()
374     plt.plot(epochs, self.losses, label="MSE")
375     plt.xlabel("Epochs")
376     plt.ylabel("MSE")
377     plt.legend(loc="upper left")
378     plt.show()
379
380 def plot_epsilon(self):
381     """Plot the total non-discounted sum of rewards across the
382     ↪ episodes (i.e duration of each episode in steps)."""
383     epochs = np.linspace(0, len(self.epsilon_list),
384     ↪ len(self.epsilon_list))
385
386     plt.figure(figsize=(10,7))
387     plt.grid()
388     plt.plot(epochs, self.epsilon_list, color="red")
389     plt.xlabel("Epochs")
390     plt.ylabel("Epsilon")
391     plt.show()
392
393 def replicate_CPAgent(N_REPLICATIONS, NUM_EPISODES, BATCH_SIZE, GAMMA,
394     ↪ TARGET_UPDATE_FREQ,
395     NUM_FRAMES, SKIPPED_FRAMES, NUM_HIDDEN_LAYERS,
396     ↪ SIZE_HIDDEN_LAYERS,
397     REPLAY_BUFFER, LR, EPS_START, EPS_END,
398     ↪ REWARD_TARGET, return_epsilon=False):
399
400     """
401     replication_rewards = []
402     replication_epsilons = []
403     for i in range(N_REPLICATIONS):
404         print(f"Replication {i+1}")
405         CP_agent = CartpoleAgent(NUM_EPISODES, BATCH_SIZE, GAMMA,
406             ↪ TARGET_UPDATE_FREQ,
407             NUM_FRAMES, SKIPPED_FRAMES,
408             ↪ NUM_HIDDEN_LAYERS,
409             ↪ SIZE_HIDDEN_LAYERS,
410             REPLAY_BUFFER, LR, EPS_START,
411             ↪ EPS_END, REWARD_TARGET,
412             ↪ random=True)

```

```

401     CP_agent.train() #train agent
402     replication_rewards.append(CP_agent.train_rewards)
403     if return_epsilon:
404         replication_epsilon_list.append(CP_agent.epsilon_list)
405
406 mean_rewards = np.mean(np.array(replication_rewards), axis=0)
407 std_rewards = np.std(np.array(replication_rewards), axis=0)
408
409 if return_epsilon:
410     mean_epsilon_list = np.mean(np.array(replication_epsilon_list),
411                                ↪ axis=0)
412     std_epsilon_list = np.std(np.array(replication_epsilon_list), axis=0)
413
414     return mean_rewards, std_rewards, mean_epsilon_list, std_epsilon_list
415 else:
416     return mean_rewards, std_rewards
417
418
419 def plot_replications(mean_rewards, std_rewards, mean_epsilon_list=None,
420                      ↪ std_epsilon_list=None):
421     """
422     episodes = np.arange(len(mean_rewards))
423
424     plt.figure(figsize=(10,7))
425     plt.grid()
426     plt.plot(episodes, mean_rewards, "r", label="mean /  $\mu$ ")
427     plt.fill_between(episodes, mean_rewards-std_rewards,
428                    ↪ mean_rewards+std_rewards, color="r", alpha=0.4, label=r"std /
429                    ↪  $\sigma$ ")
430     plt.xlabel("Episode")
431     plt.ylabel("Mean Total Non-Discounted Reward")
432     plt.ylim((0,550))
433     plt.axhline(y=500, color='black', linestyle='--', label="Maximum
434                    ↪ reward")
435     plt.legend(loc="best")
436     plt.show();
437
438 if mean_epsilon_list is not None and std_epsilon_list is not None:
439     plt.figure(figsize=(10,7))
440     plt.grid()
441     plt.plot(episodes, mean_epsilon_list, "b", label="mean /  $\mu$ ")
442     plt.fill_between(episodes, mean_epsilon_list-std_epsilon_list,
443                    ↪ mean_epsilon_list+std_epsilon_list, color="b", alpha=0.4,
444                    ↪ label=r"std /  $\sigma$ ")

```

```
439         plt.xlabel("Episode")
440         plt.ylabel("Epsilon")
441         plt.legend(loc="best")
442         plt.show();
443
444
445 if __name__ == "__main__":
446     #define hyperparameters for agent
447     NUM_EPISODES = 700
448     BATCH_SIZE = 128
449     GAMMA = 0.9999
450     TARGET_UPDATE_FREQ = 40
451     NUM_FRAMES = 4
452     SKIPPED_FRAMES = 4
453     NUM_HIDDEN_LAYERS = 2
454     SIZE_HIDDEN_LAYERS = 150
455     REPLAY_BUFFER = 100000
456     LR = 0.0001
457     EPS_START = 1.0
458     EPS_END = 0.0005
459     REWARD_TARGET = 100
460
461     CP_agent = CartpoleAgent(NUM_EPISODES, BATCH_SIZE, GAMMA,
462                               ↪ TARGET_UPDATE_FREQ,
463                               NUM_FRAMES, SKIPPED_FRAMES,
464                               ↪ NUM_HIDDEN_LAYERS,
465                               ↪ SIZE_HIDDEN_LAYERS,
466                               REPLAY_BUFFER, LR, EPS_START, EPS_END,
467                               ↪ REWARD_TARGET)
468
469     CP_agent.train()
470     CP_agent.plot_rewards()
```
