

1 Question 1: Implementing a functional DQN

1.1 Implementing DQN (See Appendix)

1.2 Implementing Atari game technical features

Firstly we implement a replay buffer to model experience replay. The replay buffer contains a collection of experience tuples (s, a, s', r) . The tuples are gradually added to the buffer as we are interacting with the Environment. We implement a fixed size buffer, with new data added to the end of the buffer so that it pushes the oldest experience out of memory. A 'deque' data structure was used from Python's built-in collections library, this is a list that you can set a maximum size on. For training, we randomly sample the batch of transitions from the replay buffer, which allows us to break any correlation between subsequent steps in the environment.

Next I created a target network, by which we keep a copy of our neural network and use it for the $Q(s', a')$ value in the Bellman equation in order to avoid instability from the neural network updating parameters at every step. This is done by first initialising the a target network along with the policy network (lines 229-238). The predicted Q values of this second Q -network are used to backpropagate through and train the main Q -network. The target network's parameters are not trained, but they are periodically synchronized with the parameters of the main Q -network, at a set amount of episodes which can be considered another hyperparameter.

The final technique used was 'stacking frames' as used in the Atari DQN implementation. This effectively takes advantages of pushing information from four previous states to the memory at every step. This means that the neural net takes in input dimension of $k*4$ features (i.e. a $4k \times 1$ vector), alternatively this could have been pushed as a $4 \times k$ matrix, but linear neural network layers cannot take vector inputs meaning we have to flatten the four state values for each frame into one vector. This increased number of features should give the neural net more information to draw an action from, for example drawing conclusions about how the total reward is increased by adjacent frames and how the pole or cart accelerates. This could be done with the gym FrameStack wrapper but I chose to write it myself, with a `k_states` deque (of state tensors) with max length `k`. This is initialised with the environment (lines 261-263) as the first state for the first two elements. Then when an action is taken, the deque is flattened into a vector (line 266) and fed to the neural net, selecting an action from the maximum Q values. If the action selected is not terminal, the next state is appended to the end of the deque and the earliest state is shifted out of `k_states` (line 285). If the states are terminal then a tensor of zeros is added to `k_states`. `k_states` is again flattened and we push flattened_k_states, action, next_flattened_k_states and reward to memory). There is more manipulation in the optimiser to ensure that the states with a final terminal states are removed from the `non_final_mask` as per the original code.

1.3 Design Decisions

We have chosen to implement a feed-forward multi-layer neural network. This is composed of two hidden linear layers of 164 neurons with a ReLu activation function on each layer. The selection of this architecture was done by attempting a few different configurations within the range of 1-5 layers and 10-200 neurons. This experimentation demonstrated that, generally, wider networks gave a more stable learning curve alongside the other parameters. The choice of this number is dependent on the value chosen for `k` which changes the input dimension by factor `k`, and so this width and depth of network was as far as I could 'push' the architecture parameters before

some over-fitting was seen with a dramatic drop off once the agent reached high rewards. In Figure 1 the basic structure of the network is shown, with a pair of 128 neuron wide layers. I

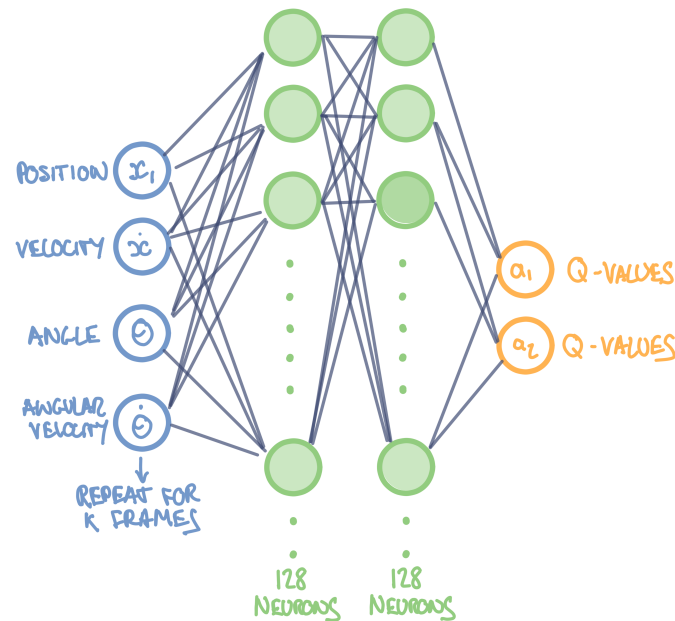


Figure 1: Neural Network Architecture

chose Adam as my optimizer. This is often the preferred optimizer used for neural nets in both classification and regression tasks using pytorch. The idea behind this method of gradient descent is to utilize the momentum concept from some Stochastic Gradient Descent implementations, and adaptive learning rate from another technique called “Ada delta”. This also appeared to give the most stable results, and particularly results that converge. I used the pytorch ‘Smooth L1 Loss’ criterion. Also known as Huber loss, it uses a squared term if the absolute error falls below 1 and an absolute term otherwise. It has the advantage over the mean squared error loss that it is less sensitive to outliers and in some cases prevents exploding gradients.

1.4 Learning Curve

Figure 2 shows the mean reward values over 10 repetitions. The hyperparameters chosen gave a steady solution that converged as well as anything tried. Whilst some runs with different parameters were able to achieve scores of 500 the agent catastrophically worsened, this can be seen in figure 3, early stopping could be implemented, but this has the problem of stopping at different points and removing the ability to replicate for the same number of episodes. In the chosen solution an average rewards of 300 was reached, with the agent stabilising after this point, and taking 125 episodes to reach 90% of the final value.

2 Question 2: Hyperparameters of the DQN

2.1 Choice of ϵ

ϵ indicates the proportion of random actions relative to taken actions that are informed by existing agent “knowledge” accumulated during the episode. This strategy is called a “Greedy Search

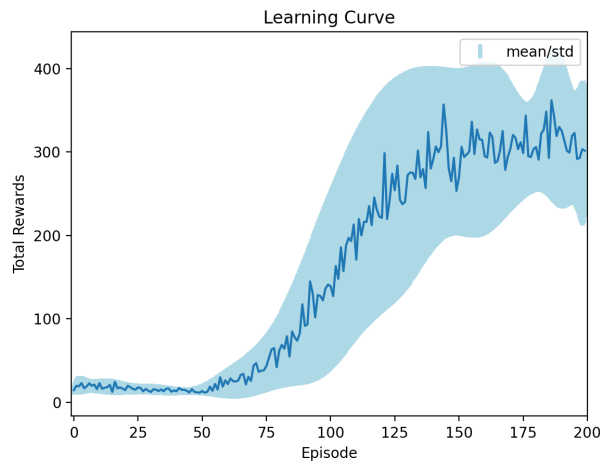


Figure 2: Learning Curve for chosen DQL parameters, showing mean and standard deviation over 10 replications for 200 episodes

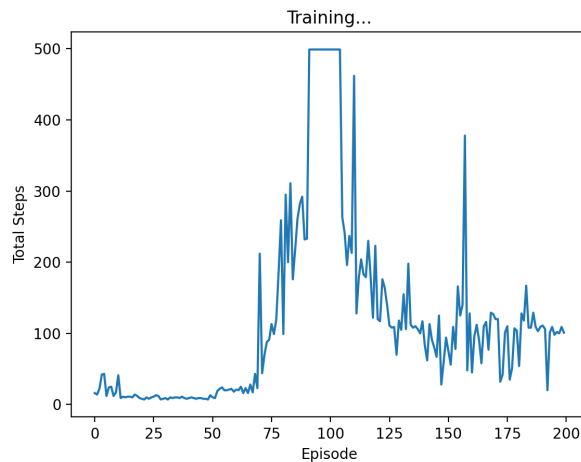


Figure 3: Learning Curve for alternative DQL parameters for one run of 200 episodes

Policy.” Before playing the game, the cart agent doesn’t have any experience, so it is common to set epsilon to higher values and then gradually decrease its value. This can decay in a number of ways, in literature usually by a constant value e.g. 0.99 every step, I have chosen a similar step but ensure that there is decay until a minimum value by creating a power law depreciation. This reduces to a value very close to zero (0.05) as more steps are taken in Figure 4 below with two other decay schedule options.

Immediately we see from the plot that for the ‘1/k’ decay, i.e. ϵ is divided by the step number in the episode, the value for ϵ reduces drastically very quickly. This will have the effect of reducing the amount of exploration the agent undertakes. The learning curve for this decay can be seen in Figure 5c). The agent can be seen to start by chance with a higher reward but then does no exploration and therefore soon has no useful knowledge. Figure 5a) shows the effect of choosing a constant ϵ value of 0.6 which was seen to work fairly well. However as epsilon doesn’t continue to decay there is still large degree of randomness after receiving a relatively high reward, the rewards fluctuate far more as it is still exploring instead of using its knowledge. The linear decay

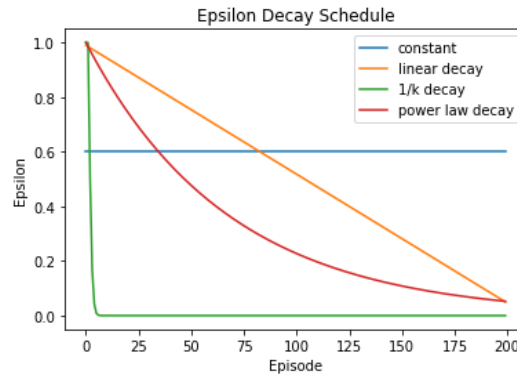


Figure 4: ϵ Decay Schedules, ϵ against N episodes

in Figure 5b) is improved but the mean fluctuates more than the power law decay, we need more exploitation sooner, and so we make a good choice for the epsilon decay schedule.

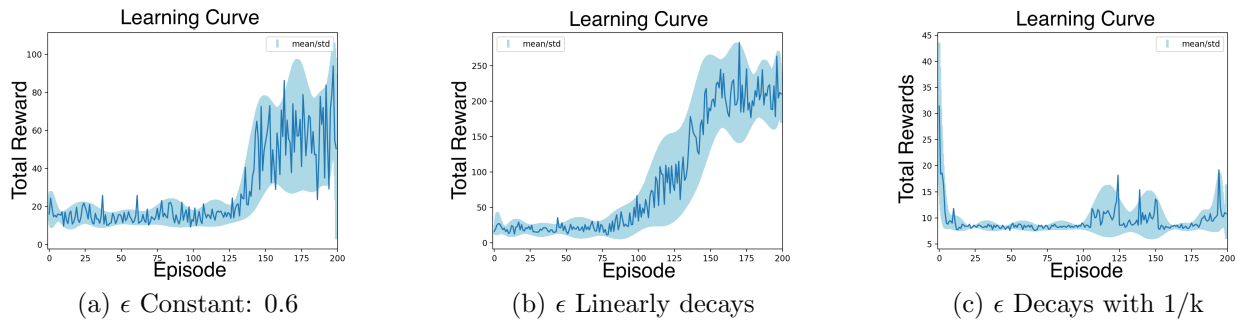


Figure 5: Variation of ϵ decay schedule on learning curve for 5 replications

2.2 Choice of Replay Buffer Size

The replay buffer gives the agent the ability to store experiences in memory. Batches of experiences are randomly sampled from memory and subsequently used to train the neural network. The size of the replay buffer is the number of experiences that are fed into the network update. It therefore makes sense that the larger the buffer, the more experiences the agent can draw from and the rewards deviate less from the maximum experienced. Figure 6 demonstrates this decrease, and particularly the large jump from 1000 to 10000 replays. The standard deviation then increases slightly for the next increment which shows that 10000 is sufficient for the agent to draw from in this context, and requires less computational resources. The standard deviation was calculated by taking 67% of the maximum total reward across all runs (for different buffer size), then working out the std with this value rather than taking the difference between each reward and the mean shown in line 480.

2.3 Choice of k

The various learning curves show that for this time period of 200 episodes, fewer frames results in quicker learning from the agent and a higher mean reward. This is shown in figure 7 below. It does however also show that for higher values of k the total rewards fluctuate less as the agent

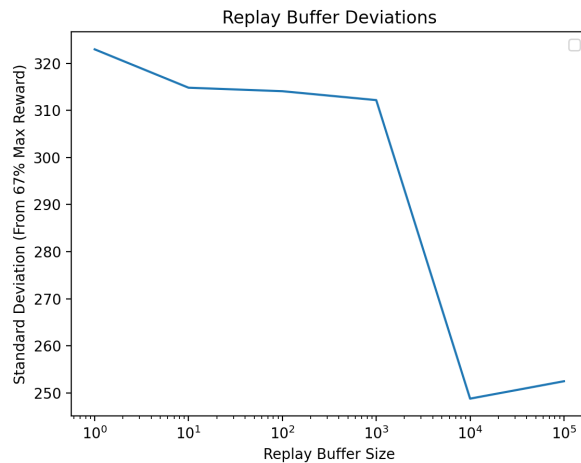


Figure 6: Standard Deviations from 67% of maximum reward (across all runs), against replay buffer size on logarithmic scale

acquires knowledge, giving a more stable learning curve. The speed of learning is expected as for lower k values the input dimension into the neural network is much smaller, with the hidden layers drawing inferences from fewer features. In practice the network architecture for each k value should be tuned to suit. In this example all values of k result in learning, and though $k=1$ appears optimal and so this value has been selected to maximise rewards. With more episodes we may have seen an improvement by using more frames, or perhaps by increasing the number of features into a flattened vector it makes the network too complex to draw intelligent features from in the neural net. It could be the case also that the CartPole environment doesn't benefit from this frame stacking as much as an Atari game using frames of pixel inputs and convolutional neural networks in comparison to our physics problem.

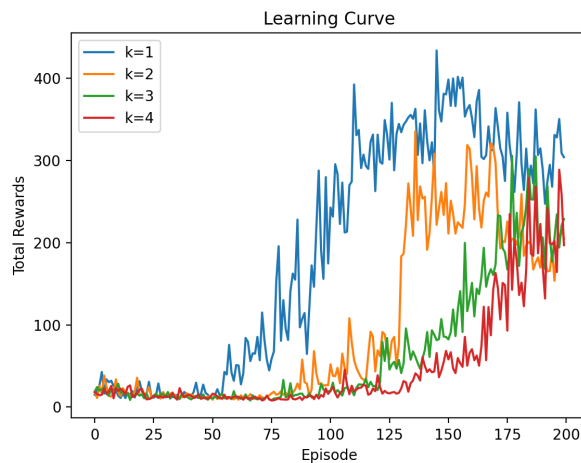


Figure 7: Variation of k on learning curve for DQL, each value of k was used for three runs and a mean value taken keeping all other hyperparameters constant

3 Question 3: Ablation/Augmentation experiments

3.1 Double Deep Q-Learning Implementation

Q-learning when introduced in a Deep Reinforcement learning context is known to overestimate the action values in certain conditions. The Double Q learning algorithm is known to reduce susceptibility to this overestimation by adjusting the error calculation using the weights of the target network to fairly evaluate the value of this policy. The implementation can be seen in lines 124-127. If the DDQL condition is True, instead of selecting the next state values from doing a forward pass of the target network and selecting the maximum value, we select the maximum action index from using the policy network weights then use the second set of weights (of the target network) to fairly evaluate the value of this policy. Note that we can switch the roles of the target and policy set of weights symmetrically according to the Deepmind paper (Van Hessel et al. 2015) where the proposition of using DDQL is advocated.

However when experimenting with this I found huge variations between implementing DDQL when the target network and policy network roles are switched. Figure 8 demonstrates this, and actually the reverse case appears to get maximum rewards more often than the action selection method specified in the paper. I chose to implement the conventional one as it still shows an improvement on the DQL learning curve with more consistency. As for all of these experiments, more episodes and some form of early stopping would vastly improve the performance of an agent but for demonstration purposes the selected hyperparameters and run conditions will suffice.

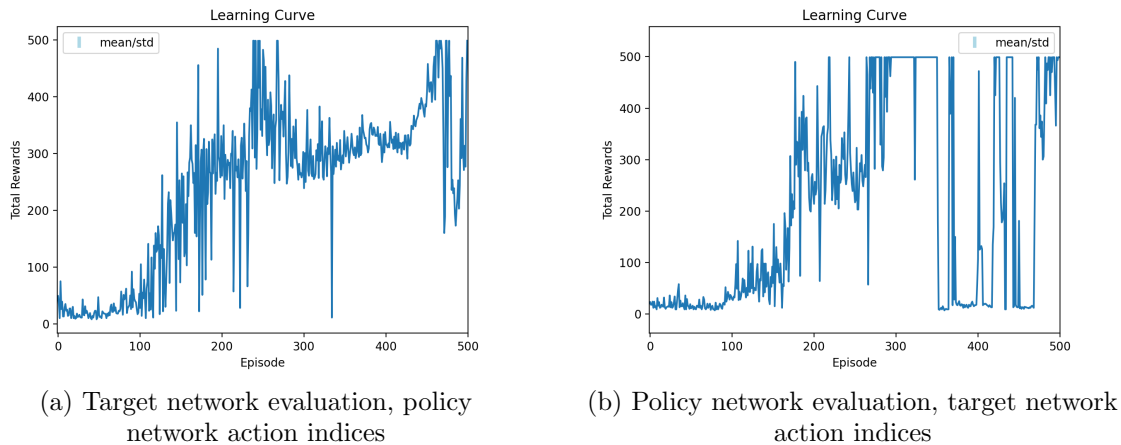


Figure 8: DDQL Implementation Reverse

3.2 Learning Curve Comparison

The chosen value of k has a similar effect as for the single DQL network, this can be seen in figure 9 with $k=1$ clearly performing best and reaching maximum reward numerous times, again fluctuation may be seen as an issue but the agent evidently maximises rewards when using only 4 input features.

Comparing the different ablation/augmentation experiments we can draw a number of conclusions from incorporating additional features to our network. Firstly, removing the replay buffer makes this a simple temporal difference learning model taking information only from the last step. This results in some learning but the speed of learning is slow. We saw a similar effect in figure 6.

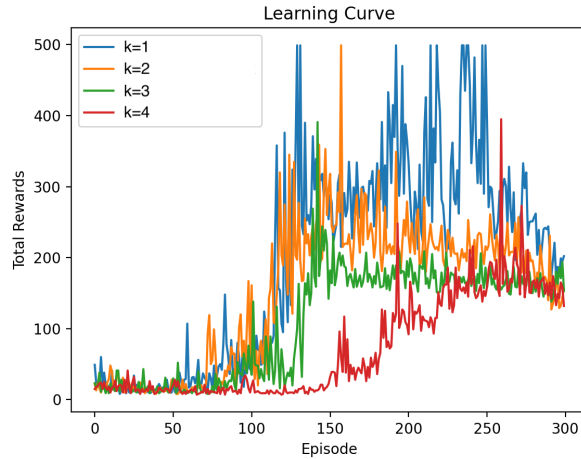


Figure 9: Variation of k on learning curve for DDQL, each value of k was used for three runs and a mean value taken keeping all other hyperparameters constant

For the other three curves we see the neural network with the replay seems far more robust and intelligent compared to its counterpart that only remembers the last action. Next we see removing the target network feature, removes the 'resonance damper' effect on the the normal curve, we see more fluctuations beyond the 150 episode mark, and while the curve may reach some of the maximum rewards experienced with the target network enabled, the curve doesn't converge with the normal or DDQN curves. This relaxation time can be seen to improve the performance and learning speed of the network. Implementing the DDQN can be seen to have positive effects on the learning of the agent, both in terms of speed and maximum reward. For a larger number of episodes as seen in figure 8a), we can see that the overall performance seems to increase to a high reward, hitting the maximum multiple times. We therefore can conclude that this method of ensuring no over-reliance on one network results on in a better network using DDQL.

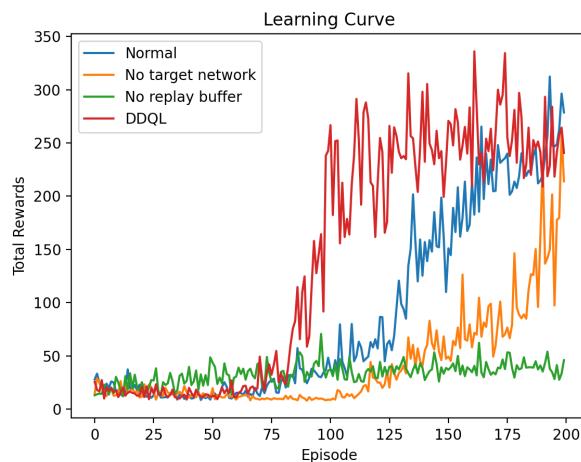


Figure 10: Learning curves incorporating ablation of target network, replay buffer and augmentation of DDQL against original DQL learning curve; using 3 runs and taking a mean score for each and $k=1$

A Appendix - Code

Notes on running the code:

- The code has been set up to run various experiments and therefore has been refactored so that there are only 4 lines of code to run at the bottom of the script (641 for a normal DQL learning curve, 650 to run the replay buffer experiment, 654 to run the k frames experiment and 658 to run the ablation/augmentation experiment). These will need to be uncommented to run, other than run_learning_curve which is left active for 10 runs
- Hyperparameters are actually set in the functions themselves which is not the best practice but convenient for running these experiments
- One input to the training loop is the 'record' toggle, which I have set to True by default now. Note that the error will be raised that video/video.mp4 filepath does not exist, this will need to be adjusted if you want to see the video to a valid filepath.
- There also exists separate plot functions, the inputs of which should be intuitive.

```
1 from IPython.display import clear_output
2 import gym
3 from gym.wrappers.monitoring.video_recorder import VideoRecorder # records videos of episodes
4 import numpy as np
5 import matplotlib.pyplot as plt # Graphical library
6 import torch
7 import torch.optim as optim
8 import torch.nn as nn
9 import torch.nn.functional as F
10
11 device = torch.device("cuda" if torch.cuda.is_available() else "cpu") # Configuring Pytorch
12 from collections import namedtuple, deque
13 from itertools import count
14 import random
15
16 clear_output()
17
18 random.seed(10)
19 Transition = namedtuple('Transition',
20                        ('state', 'action', 'next_state', 'reward'))
21
22
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
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, hidden_size) \
45                                             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
58 def optimize_model(memory,
59                   BATCH_SIZE,
60                   state_dim,
61                   policy_net,
62                   target_net,
63                   GAMMA,
64                   optimizer,
65                   DDQL,
66                   target_net_on):
67     if len(memory) < BATCH_SIZE:
68         return
69     transitions = memory.sample(BATCH_SIZE)
70     # Transpose the batch (see https://stackoverflow.com/a/19343/3343043 for
71     # detailed explanation). This converts batch-array of Transitions
72     # to Transition of batch-arrays.
73     batch = Transition(*zip(*transitions))
74
75     # Compute a mask of non-final states and concatenate the batch elements

```

```

76  # (a final state would've been the one after which simulation ended)
77  # Alteration from original code to account for k frame stacking, now
78  # treats a state as final if the last 4 elements are equal to zero.
79  non_final_mask = torch.tensor(tuple(map(lambda s: torch.sum(s[0][len(s[0]) - 4:]))
80                                  .absolute().item() > 0, batch.next_state)),
81                                  device=device,
82                                  dtype=torch.bool)
83
84  # Can safely omit the condition below to check that not all states in the
85  # sampled batch are terminal whenever the batch size is reasonable and
86  # there is virtually no chance that all states in the sampled batch are
87  # terminal. Similar change to account for k frame stacking final states
88  if sum(non_final_mask) > 0:
89      non_final_next_states = torch.cat([s for s in batch.next_state if torch.sum(
90          s[0][len(s[0]) - 4:]).absolute().item() > 0])
91  else:
92      non_final_next_states = torch.empty(0, state_dim).to(device)
93
94  state_batch = torch.cat(batch.state)
95  action_batch = torch.cat(batch.action)
96  reward_batch = torch.cat(batch.reward)
97
98  # Compute Q(s_t, a) - the model computes Q(s_t), then we select the
99  # columns of actions taken. These are the actions which would've been taken
100  # for each batch state according to policy_net
101  state_action_values = policy_net(state_batch).gather(1, action_batch)
102
103  actions_test = policy_net(state_batch)
104  # Compute V(s_{t+1}) for all next states.
105  # This is merged based on the mask, such that we'll have either the expected
106  # state value or 0 in case the state was final.
107  next_state_values = torch.zeros(BATCH_SIZE, device=device)
108
109  with torch.no_grad():
110      # Once again can omit the condition if batch size is large enough
111      if sum(non_final_mask) > 0:
112          if not DDQL:
113              # Single DQL case, with additional condition for using target
114              # network updates (for final ablation question)
115              if target_net_on:
116                  next_state_values[non_final_mask] = target_net(
117                      non_final_next_states).max(1)[0].detach()
118              else:
119                  next_state_values[non_final_mask] = policy_net(
120                      non_final_next_states).max(1)[0].detach()
121      else:

```

```

122         # DDQL case, select action from policy network and then values from
123         # target network
124         policy_actions = \
125             policy_net(non_final_next_states).max(1)[1].view(-1, 1)
126         next_state_values[non_final_mask] = \
127             target_net(non_final_next_states).gather(1, policy_actions).view(-1, )
128     else:
129         next_state_values = torch.zeros_like(next_state_values)
130
131     # Compute the expected Q values
132     expected_state_action_values = (next_state_values * GAMMA) + reward_batch
133
134     # Compute loss using mean squared error loss criterion CHANGE IN REPORT
135     criterion = nn.MSELoss()
136     loss = criterion(state_action_values, expected_state_action_values.unsqueeze(1))
137
138     # Optimize the model
139     optimizer.zero_grad()
140     loss.backward()
141
142     # Limit magnitude of gradient for update step (THIS LIMITATION HAS BEEN REMOVED,
143     # SLOWING DOWN RUN TIME BUT IMPROVING RESULTS)
144     # for param in policy_net.parameters():
145     #     param.grad.data.clamp_(-1, 1)
146
147     optimizer.step()
148
149
150 def plot_total_rewards(N_episodes, total_rewards):
151     plt.figure(2)
152     episodes = np.arange(N_episodes)
153     plt.title('Training...')
154     plt.xlabel('Episode')
155     plt.ylabel('Total Steps')
156     plt.plot(episodes, np.array(total_rewards))
157     print(f"Average Reward: {sum(total_rewards) / N_episodes}")
158
159
160 def select_action(state=None, current_eps=0, n_actions=2, policy_net=None):
161     sample = random.random()
162     eps_threshold = current_eps
163     if sample > eps_threshold:
164         with torch.no_grad():
165             # t.max(1) will return largest column value of each row.
166             # second column on max result is index of where max element was
167             # found, so we pick action with the larger expected reward.
168             return policy_net(state).max(1)[1].view(1, 1)

```

```

169     else:
170         return torch.tensor([[random.randrange(n_actions)]],
171                               device=device,
172                               dtype=torch.long)
173
174
175 def eps_decay(NUM_EPISODES, EPS_START, EPS_END, decay_type):
176     '''
177     Function to produce different epsilon decay schedules
178     '''
179
180     eps_linear = np.linspace(EPS_START, EPS_END, NUM_EPISODES)
181     eps_const = 0.6 * np.ones(eps_linear.shape)
182     eps_glie = np.ones(eps_linear.shape)
183     eps_factor = np.ones(eps_linear.shape)
184     factor = np.power(EPS_END, EPS_START / NUM_EPISODES)
185
186     for i in range(1, len(eps_linear)):
187         eps_glie[i] = eps_glie[i - 1] / i
188         eps_factor[i] = eps_factor[i - 1] * factor
189
190     if decay_type == 'linear':
191         return eps_linear
192     elif decay_type == '1/k':
193         return eps_glie
194     elif decay_type == 'const':
195         return eps_const
196     else:
197         return eps_factor
198
199
200 def train(NUM_EPISODES=100,
201           BATCH_SIZE=128,
202           GAMMA=0.99,
203           EPS_START=0.9,
204           EPS_END=0.05,
205           EPS_DECAY='power',
206           LR=0.0001,
207           num_hidden_layers=2,
208           size_hidden_layers=128,
209           network_sync_freq=10,
210           k=1,
211           replays=10000,
212           target_net_on=True,
213           DDQL=False,
214           record=False):
215     '''

```

```

216 Main training loop for agent,
217 input: all hyperparameters and ablation/augmentation toggles
218 return: (np array) total rewards for each episode
219 '''
220
221 # Get number of states and actions from gym action space
222 env = gym.make("CartPole-v1")
223 env.reset()
224 state_dim = k * len(env.state) # x, x_dot, theta, theta_dot
225 n_actions = env.action_space.n
226 env.close()
227
228 # Initialise two networks, policy net and identical target net
229 policy_net = DQN(state_dim,
230                 n_actions,
231                 num_hidden_layers,
232                 size_hidden_layers).to(device)
233 target_net = DQN(state_dim,
234                 n_actions,
235                 num_hidden_layers,
236                 size_hidden_layers).to(device)
237 target_net.load_state_dict(policy_net.state_dict())
238 target_net.eval()
239
240 optimizer = optim.Adam(policy_net.parameters(), LR)
241 memory = ReplayBuffer(replays)
242
243 # Empty list to append total rewards for each episode to
244 # (same as duration of episode)
245 durations = []
246 epsilon = EPS_START
247
248 # Use custom function to get decay schedule
249 eps_schedule = eps_decay(NUM_EPISODES, EPS_START, EPS_END, EPS_DECAY)
250
251 for i_episode in range(NUM_EPISODES):
252
253     if i_episode % 20 == 0:
254         print("episode ", i_episode, "/", NUM_EPISODES)
255
256     # Initialize the environment and state
257     env.reset()
258     state = torch.tensor(env.state).float().unsqueeze(0).to(device)
259
260     # Initialise frame stacker and set first k frames as initial state
261     k_states = deque([], maxlen=k)
262     for i in range(k):

```

```
263     k_states.append(state)
264
265     for t in count():
266         flattened_k_states = torch.stack(
267             list(k_states)).reshape(-1).unsqueeze(0) # Flatten
268         action = select_action(flattened_k_states,
269                               epsilon,
270                               n_actions,
271                               policy_net) # Select action from network
272         _, reward, done, _ = env.step(action.item())
273         reward = torch.tensor([reward], device=device)
274
275         # Observe new state
276         if not done:
277             next_state = torch.tensor(
278                 env.state).float().unsqueeze(0).to(device)
279         else:
280             # If terminal set next state as zeros
281             next_state = torch.zeros_like(state)
282
283         # Store the transition in memory
284         # Append state to frame stacker deque, pushing out oldest frame
285         k_states.append(next_state)
286         next_flattened_k_states = torch.stack(
287             list(k_states)).reshape(-1).unsqueeze(0) # Flatten k states again
288         memory.push(flattened_k_states,
289                   action,
290                   next_flattened_k_states,
291                   reward) # Push (s,a,s',r) to memory
292
293         # Move to the next state
294         state = next_state
295
296         # Perform one step of the optimization (on the policy network)
297         optimize_model(memory,
298                       BATCH_SIZE,
299                       state_dim,
300                       policy_net,
301                       target_net,
302                       GAMMA,
303                       optimizer,
304                       DDQL,
305                       target_net_on)
306     if done:
307         break
308
```

```
309     # Move onto next epsilon value in schedule
310     epsilon = eps_schedule[i_episode]
311
312     durations.append(t)
313
314     # Sync target network with policy net every set number of episodes
315     if i_episode % network_sync_freq == 0:
316         target_net.load_state_dict(policy_net.state_dict())
317
318     print('Complete')
319
320     env.close()
321
322     # Plot rewards against episodes
323     # plot_total_rewards(NUM_EPISODES, durations)
324     # plt.show()
325
326     # If record toggle on, record the episode
327     if record:
328         record_cart(policy_net, k)
329
330     return durations
331
332
333 def record_cart(policy_net, k):
334     '''
335     Record function separated to be toggled in main training loop, if you want
336     to see the video.
337     '''
338     env = gym.make("CartPole-v1")
339     file_path = 'video/video.mp4'
340     recorder = VideoRecorder(env, file_path)
341
342     observation = env.reset()
343     done = False
344
345     state = torch.tensor(env.state).float().unsqueeze(0)
346     k_states = deque([], maxlen=k)
347     for i in range(k):
348         k_states.append(state)
349
350     duration = 0
351
352     while not done:
353         recorder.capture_frame()
354
355         # Select and perform an action
```

```

356     flattened_k_states = torch.stack(
357         list(k_states)).reshape(-1).unsqueeze(0)
358     action = select_action(flattened_k_states,
359                           current_eps=0,
360                           n_actions=2,
361                           policy_net=policy_net)
362
363     observation, reward, done, _ = env.step(action.item())
364     duration += 1
365     reward = torch.tensor([reward], device=device)
366
367     # Observe new state
368     state = torch.tensor(env.state).float().unsqueeze(0)
369     k_states.append(state)
370
371     recorder.close()
372     env.close()
373     print("Episode duration: ", duration)
374
375
376 def plot_learning_rate_full(total_rewards):
377     '''
378     Function to plot learning curves for a number of different training runs,
379     taking mean and standard deviation and plotting this on the figure as a
380     line with errorbars.
381     Total rewards given as a numpy array (R x N) with each row as a full
382     training run under certain conditions:
383     R = number of runs
384     N = number of episodes per run
385     '''
386     n_episodes = np.shape(total_rewards)[1]
387     episodes = np.arange(0, n_episodes, 1)
388
389     # Take mean and std across different runs
390     mean = np.mean(total_rewards, axis=0)
391     std = np.std(total_rewards, axis=0)
392
393     # MEAN
394     x = episodes
395     y1 = mean
396     y2 = std
397
398     # calculate polynomial
399     z1 = np.polyfit(x, y1, 30)
400     z2 = np.polyfit(x, y2, 30)
401     f1 = np.poly1d(z1)

```



```

402     f2 = np.poly1d(z2)
403
404     # calculate new x's and y's
405     x_new = np.linspace(x[0], x[-1], 500)
406     y_new1 = f1(x_new)
407     y_new2 = f2(x_new)
408
409     plt.plot(x, mean, '-', markersize=1)
410     plt.errorbar(x_new,
411                  y_new1,
412                  yerr=y_new2,
413                  fmt='none',
414                  color='blue',
415                  ecolor='lightblue',
416                  elinewidth=3,
417                  capsize=0,
418                  label="mean/std")
419
420     plt.xlim([x[0] - 1, x[-1] + 1])
421     plt.xlabel('Episode')
422     plt.ylabel("Total Rewards")
423     plt.legend()
424     plt.title("Learning Curve")
425     plt.show()
426
427
428 def plot_k_learning_rates(k_total_rewards):
429     '''
430     Function to plot mean learning curves for different values of k.
431     k_total_rewards given as a numpy array (R x N) with each row as a full
432     training run under certain conditions:
433     R = number of runs
434     N = number of episodes per run
435     Also used for ablation and augmentation experiment
436     '''
437
438     # Select legend labels
439     k_values = ["k=1", "k=2", "k=3", "k=4"]
440     alterations = ["Normal", "No target network", "No replay buffer", "DDQL"]
441
442     for i, total_rewards in enumerate(k_total_rewards):
443         n_episodes = np.shape(total_rewards)[1]
444         episodes = np.arange(0, n_episodes, 1)
445
446         mean = np.mean(total_rewards, axis=0)
447

```

```
448     x = episodes
449     y = mean
450
451     plt.plot(x, y, label=f"{alterations[i]}")
452
453     plt.xlabel('Episode')
454     plt.ylabel("Total Rewards")
455     plt.legend()
456     plt.title("Learning Curve")
457     plt.show()
458
459
460 def plot_replay_deviation(replays_experiment, total_rewards):
461     '''
462     Function to plot standard deviation against size of replay buffer using a
463     log x axis.
464     total_rewards: given as a numpy array (R x N) with each row as a full
465     training run under certain conditions
466     R = number of runs
467     N = number of episodes per run
468
469     replays_experiment: given as list of replay buffer sizes
470     '''
471
472     n_episodes = np.shape(total_rewards)[1]
473     n_runs = np.shape(total_rewards)[0]
474     max_reward = np.max(total_rewards)
475     level = 0.67 * max_reward
476     stds = []
477
478     for run in range(n_runs):
479         rewards = total_rewards[run, :]
480         std = np.sqrt(np.sum((rewards - level) ** 2) / n_episodes)
481         stds.append(std)
482
483     x = replays_experiment
484     y = np.array(stds)
485
486     plt.plot(x, y)
487
488     plt.xlabel('Replay Buffer Size')
489     plt.xscale('log')
490     plt.ylabel("Standard Deviation (From 67% Max Reward)")
491     plt.legend()
492     plt.title("Replay Buffer Deviations")
493     plt.show()
494
```

```
495
496 # LEARNING CURVE
497 def run_learning_curve(NUM_EPISODES, N_RUNS, k, replays, target, DDQL):
498     '''
499     Hyperparameters selected here. Train agent over specified number of N_RUNS
500     and plot the averaged learning curve return total_rewards as an (R x N)
501     matrix:
502     R = number of runs
503     N = number of episodes per run
504     '''
505     total_rewards = np.zeros(NUM_EPISODES)
506
507     for run in range(N_RUNS):
508         print(f"Training run: {run}")
509         rewards = train(NUM_EPISODES,
510                         BATCH_SIZE=32,
511                         GAMMA=0.99,
512                         EPS_START=0.99,
513                         EPS_END=0.05,
514                         EPS_DECAY='power',
515                         LR=0.0001,
516                         num_hidden_layers=2,
517                         size_hidden_layers=128,
518                         network_sync_freq=10,
519                         k=k,
520                         replays=replays,
521                         target_net_on=target,
522                         DDQL=DDQL,
523                         record=True)
524
525         total_rewards = np.vstack((total_rewards, np.array(rewards)))
526
527     total_rewards = np.delete(total_rewards, (0), axis=0)
528     plot_learning_rate_full(total_rewards)
529
530     return total_rewards
531
532
533 # REPLAYS EXPERIMENT
534 def run_replays_experiment(NUM_EPISODES, replays_experiment):
535     '''
536     Train agent over different replay buffer sizes as a list
537     (replays_experiment), and plot results. Hyperparameters selected here.
538     '''
539     total_rewards = np.zeros(NUM_EPISODES)
540
541     for run in range(len(replays_experiment)):
```

```

542     replays = replays_experiment[run]
543     print(f"Replay Buffer: {replays}")
544     rewards = train(NUM_EPISODES=200,
545                     BATCH_SIZE=32,
546                     GAMMA=0.99,
547                     EPS_START=0.99,
548                     EPS_END=0.05,
549                     EPS_DECAY='power',
550                     LR=0.0001,
551                     num_hidden_layers=2,
552                     size_hidden_layers=128,
553                     network_sync_freq=10,
554                     k=3,
555                     replays=replays,
556                     DDQL=False,
557                     record=False)
558
559     total_rewards = np.vstack((total_rewards, np.array(rewards)))
560
561     total_rewards = np.delete(total_rewards, (0), axis=0)
562     plot_replay_deviation(replays_experiment, total_rewards)
563
564
565 # K EXPERIMENT
566 def run_k_experiment(NUM_EPISODES, N_RUNS, k_experiment):
567     '''
568     Train agent over different k values (number of frames to stack) as a list
569     (k_experiment), run each for N_RUNS and plot mean results.
570     '''
571     k_total_rewards = []
572     total_rewards = np.zeros(NUM_EPISODES)
573
574     for run in range(len(k_experiment)):
575         k = k_experiment[run]
576         print(f"k value: {k}")
577
578         total_rewards = run_learning_curve(NUM_EPISODES,
579                                           N_RUNS,
580                                           k,
581                                           replays=10000,
582                                           target=True,
583                                           DDQL=True)
584         k_total_rewards.append(total_rewards)
585
586     plot_k_learning_rates(k_total_rewards)
587

```

```
588
589 # ABLATION/AUGMENTATION EXPERIMENT
590 def run_ab_experiment(NUM_EPISODES, N_RUNS):
591     '''
592     Train agent ablating and augmenting different features. Run each for N_RUNS
593     and plot all on same graph.
594     '''
595     altered_total_rewards = []
596     total_rewards = np.zeros(NUM_EPISODES)
597
598     for run in range(4):
599         if run == 0:
600             print("Normal run")
601             total_rewards = run_learning_curve(NUM_EPISODES,
602                                                 N_RUNS,
603                                                 k=1,
604                                                 replays=10000,
605                                                 target=True,
606                                                 DDQL=False)
607         elif run == 1:
608             print("Ablating target network feature")
609             total_rewards = run_learning_curve(NUM_EPISODES,
610                                                 N_RUNS,
611                                                 k=1,
612                                                 replays=10000,
613                                                 target=False,
614                                                 DDQL=False)
615         elif run == 2:
616             print("Ablating replay buffer")
617             total_rewards = run_learning_curve(NUM_EPISODES,
618                                                 N_RUNS,
619                                                 k=1,
620                                                 replays=1,
621                                                 target=True,
622                                                 DDQL=False)
623         elif run == 3:
624             print("Implementing DDQN")
625             total_rewards = run_learning_curve(NUM_EPISODES,
626                                                 N_RUNS,
627                                                 k=1,
628                                                 replays=10000,
629                                                 target=True,
630                                                 DDQL=True)
631
632         altered_total_rewards.append(total_rewards)
633
634     plot_k_learning_rates(altered_total_rewards)
```

```
635
636
637 # Run various 'experiments' in main run code below. Uncomment experiments that
638 # are not needed.
639
640 # Q1 IMPLEMENT DQN SOLUTION
641 run_learning_curve(NUM_EPISODES=200,
642                   N_RUNS=10,
643                   k=1,
644                   replays=10000,
645                   target=True,
646                   DDQL=False)
647
648 # Q2 Hyperparameters of the DQN
649 # Run replays experiment for chosen parameters
650 # run_replays_experiment(NUM_EPISODES=200,
651 #                        replays_experiment = [1, 10, 100, 1000, 10000, 100000])
652
653 # Run k experiment for chosen parameters
654 # run_k_experiment(NUM_EPISODES=300, N_RUNS=1, k_experiment=[1, 2, 3, 4])
655
656 # Q3 Ablation/Augmentation Experiments
657 # Run ablation/augmentation experiment for chosen parameters
658 # run_ab_experiment(NUM_EPISODES=200, N_RUNS=3)
```