

Reinforcement Learning

DQN for Pong – Project Report

Elise Chin, Sanjith Bonela

May 23, 2021

1 Deep Q-learning method

The deep Q-learning method was first introduced by Deepmind in 2013 [2] and then further improved and elaborated in 2015 [1]. The algorithm is an enhancement of the classic RL algorithm called Q-learning with neural networks used as a function approximator of the optimal state-action value function.

1.1 Q-Learning

Q-learning is an algorithm to find an estimate Q of the optimal state-action value function q_* from experience. The q -function of a policy π , $q_\pi(s, a)$, corresponds to the expected return or discounted sum of rewards of taking the action a at state s first, and then following the policy π . In other words, it tells us how good it is to take this specific action at that particular state. The optimal q -function is then defined as the maximum return that can be obtained starting from observation s , taking action a , and then following the optimal policy π_* . It obeys the Bellman optimality equation:

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a]$$

where R_{t+1} is the immediate reward, S_{t+1} the next state after taking action a from state s at time t .

The Q-learning update is based on the Bellman optimality equation and is defined as:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$$

and it can be shown that this converges to the optimal action-value function $q_*(s, a)$ as $N(s, a) \rightarrow \infty$ if the step size α decreases towards 0 with a suitable rate.

1.2 Deep Q-Learning

In Q-learning, a memory table $Q[s, a]$ is built to store all q -values for all possible combinations of states and actions. However, if the state or action space is too large, the memory and the computation requirement for Q-learning will be very high. To address this problem, one can train a function approximator, such as a neural network with parameters \mathbf{w} , to estimate the q -values, i.e. $Q(s, a, \mathbf{w}) \approx q_*(s, a)$. This can be done by minimizing the mean squared error at each step t :

$$L(\mathbf{w}_t) = \mathbb{E}[(U_t - Q(S_t, A_t, \mathbf{w}_t))^2]$$

where $U_t = R_{t+1} + \gamma \max_a Q(S_{t+1}, a, \mathbf{w}_t)$ is the TD-target. The loss function can be minimized using Stochastic (semi-)gradient descent:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha[U_t - Q(S_t, A_t, \mathbf{w}_t)]\nabla Q(S_t, A_t, \mathbf{w}_t)$$

However, the training will be unstable if we apply the standard Q-learning update naively. As we can see in the formula, there are some correlation between successive updates, and the estimate Q and the target are not completely independent from each other.

To address that, two ideas were introduced in the original paper [1].

1. Experience replay: store the agent's experience at each time step t , $e_t = (S_t, A_t, R_t, S_{t+1})$, in a replay buffer and sample mini-batches from it to train the network, instead of only using the last transition to compute the loss and its gradient. This has two advantages: better stability since the transitions in a batch are uncorrelated, and better data efficiency since an experience can be used in many weight updates.
2. Target network: use a separate network to estimate the TD-target. The target network has the same architecture as the function approximator but with frozen weights for a number of steps. More precisely, every C steps (a hyperparameter), the parameters of the Q-network is copied to the target network, and the latter is used to compute the TD-target for C steps. This adds more stability to the training since the target function is fixed for a certain time.

2 DQN for CartPole-v0

Since DQN can be tricky to debug, we first implement it to solve a simple environment, the CartPole-v0.

2.1 Description of CartPole-v0

A pole is attached to a cart, which can move along a frictionless track. The pole starts upright and the goal is to prevent it from falling over by controlling the cart.¹

- The observation from the environment is a 4D vector representing the position and velocity of the cart, and the angle and angular velocity of the pole.
- The agent can control the system by taking one of 2 actions: push the cart right or left.
- A reward of 1 is provided for every timestep that the pole remains upright.
- The episode ends when one of the following is true:
 - the pole tips over some angle limit
 - the cart moves outside of the world edges
 - 200 time steps pass

The goal of the agent is to find a policy as to maximize the expected discounted return $G = \sum_{t=0}^T \gamma^t R_{t+1}$ of each episode, where $\gamma \in (0, 1]$ is the discounted factor. If the agent has been trained well, an episode should end by satisfying the third condition. Hence, the final agent should have a mean return of +200.0 over several episodes, which corresponds to receiving a reward of 1 for 200 time steps.

2.2 Results with default hyperparameters

We use the following hyperparameters to train the DQN for 1000 episodes.

¹Description from https://www.tensorflow.org/agents/tutorials/0_intro_rl, <https://gym.openai.com/envs/CartPole-v0/>

Hyperparameter	Value	Description
Memory size	50000	Maximum size of the replay buffer.
Batch size	32	Size of mini-batches sampled from the replay buffer.
Target update frequency	100	The frequency (measured in time steps) with which the target network is updated.
Train frequency	1	The frequency (measured in time steps) with which the network is trained.
Gamma	0.95	Discounted factor gamma used in the Q-learning update.
Learning rate	0.0001	The learning rate used by Adam optimizer.
Epsilon start	1	Initial value of ϵ -greedy exploration.
Epsilon end	0.05	Final value of ϵ -greedy exploration.
Anneal length	10^4	Value to compute the number of steps to anneal epsilon for $\text{eps_step} = (\text{eps_start} - \text{eps_end}) / \text{anneal_length}$

Below, we plot episode loss and undiscounted return over the course of training. After increasing a bit, the loss slowly starts to decrease after 200 episodes, suggesting that the network is learning. The loss seems to continue to decrease if we ran for a greater number of episodes but it was not necessary to do so as we can see on the return plot that the agent has already reached a mean return of 200 after 400 episodes.

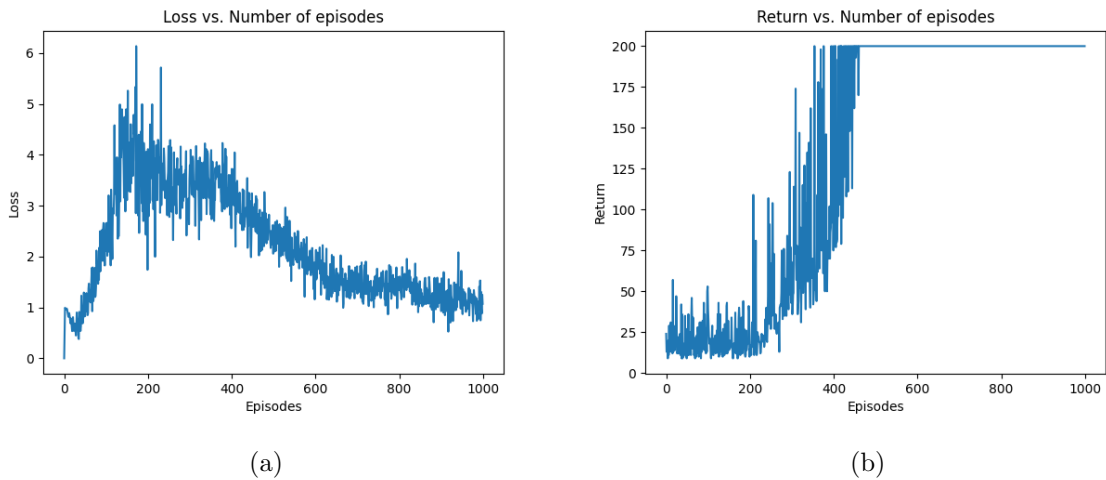


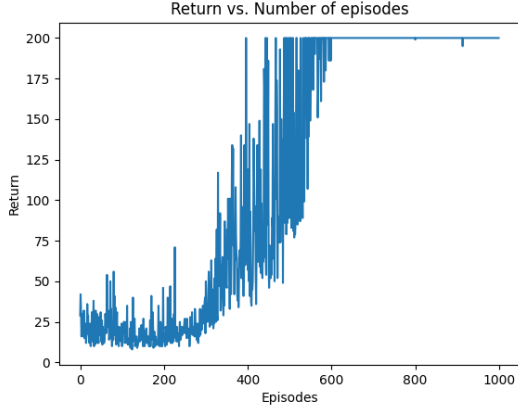
Figure 1: (a) The loss vs. the number of episodes. The network is learning well since the loss is decreasing. (b) The undiscounted return vs. the number of episodes. During the first episodes, the return is low since the agent is starting to learn, then it is increasing until it reached the maximum return which can be obtained in the CartPole-v0 environment.

2.3 Varying the target update frequency

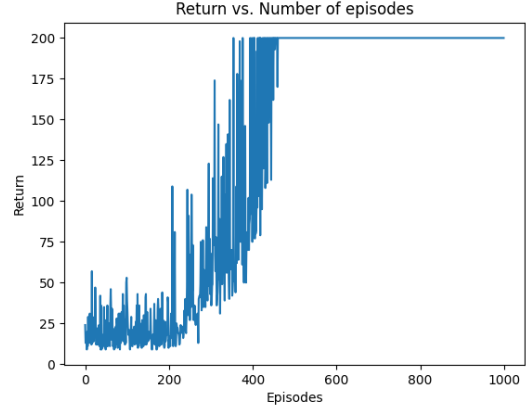
In the CartPole-v0 environment, the reward provided at each time step is 1, so the undiscounted return for an episode also corresponds to the number of steps of that episode. With a target update frequency of 100, we update the target network in the middle of an episode in average. What happens if we set this hyperparameter to a low or large value?

Below, we plot the undiscounted return over 1000 episodes, with different target update frequency (10, 100, 200, and 1000). Firstly, we can observe that the final agent always reach a return of 200 after 400 episodes regardless of the target update frequency. Secondly, the larger the target frequency is, the faster the agent learns. As we mentioned earlier in Section 2, using a separate network to compute the TD-target adds more stability to the learning. This observation can be seen in our plots where stability is reflected in the degree of fluctuation of the curves. When the target update frequency is small, the

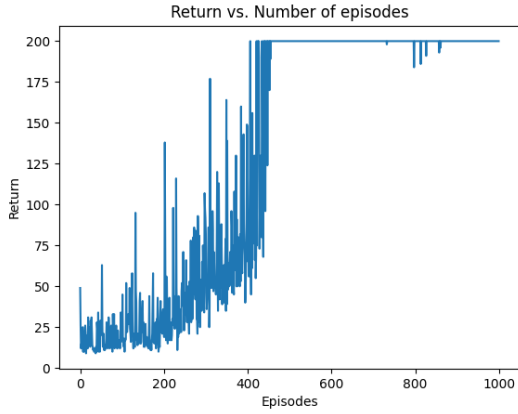
return's curve fluctuates a lot before 400 episodes. The training is more stable when the target update frequency is large because the target network is fixed for a longer period of time.



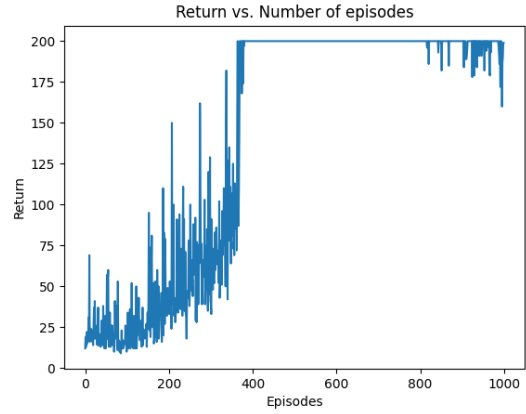
(a) target_update_frequency = 10



(b) target_update_frequency = 100



(c) target_update_frequency = 200

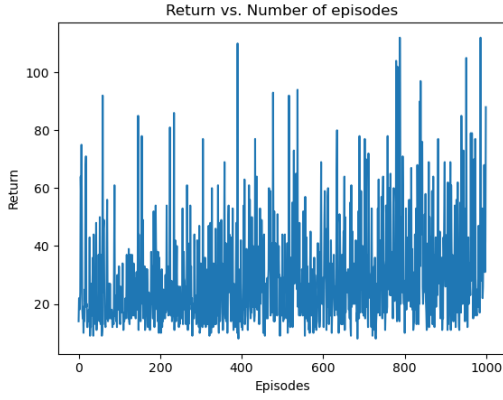


(d) target_update_frequency = 1000

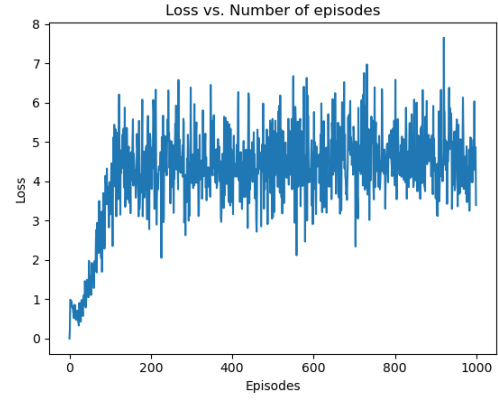
Figure 2: Undiscounted return over 1000 episodes, with different target update frequency.

2.4 Exploring different exploration parameters

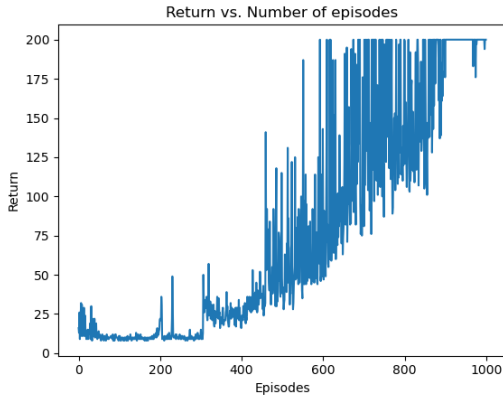
We had also explored the results with varying different hyperparameters. This exploration was concentrated around change of epsilon by varying the anneal length (which varies the decay in epsilon) in one case and exploration ($\epsilon = 1$, picks random action) and exploitation ($\epsilon = 0.05$, picks action based on max values of q) concepts in another case. The graphs are plotted below in figure [3].



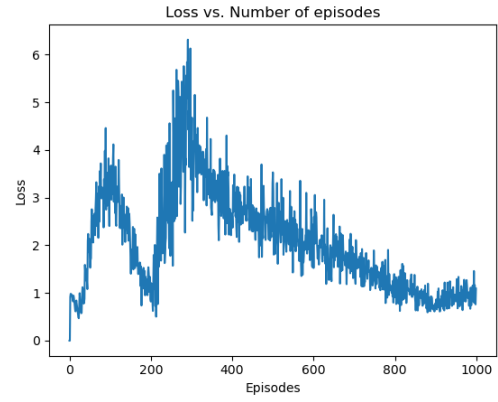
(a) Return vs. Number of episodes with anneal.length 1e-5



(b) Loss vs. Number of episodes with anneal.length 1e-5



(c) Return vs. Number of episodes with anneal.length 1e-3



(d) Loss vs. Number of episodes with anneal.length 1e-3

Figure 3: Plots for losses and returns vs Number of episodes with varying anneal lengths

From the plots in figure [3], we can see for very low anneal length ($1e-5$) (fig 3a, 3b), the loss was almost constant after around 100 episodes and the mean return was not as high as expected for 1000 episodes which may require some more training. On the contrary, if the anneal length was high ($1e-3$), the loss was constantly decreasing after around 300 episodes and the loss-episode behaviour prior to 300 episodes may account in random-action pick policy (exploration) instead of picking action that maximizes q-values and the return obtained was 200 (maximum) at around 900 episodes.

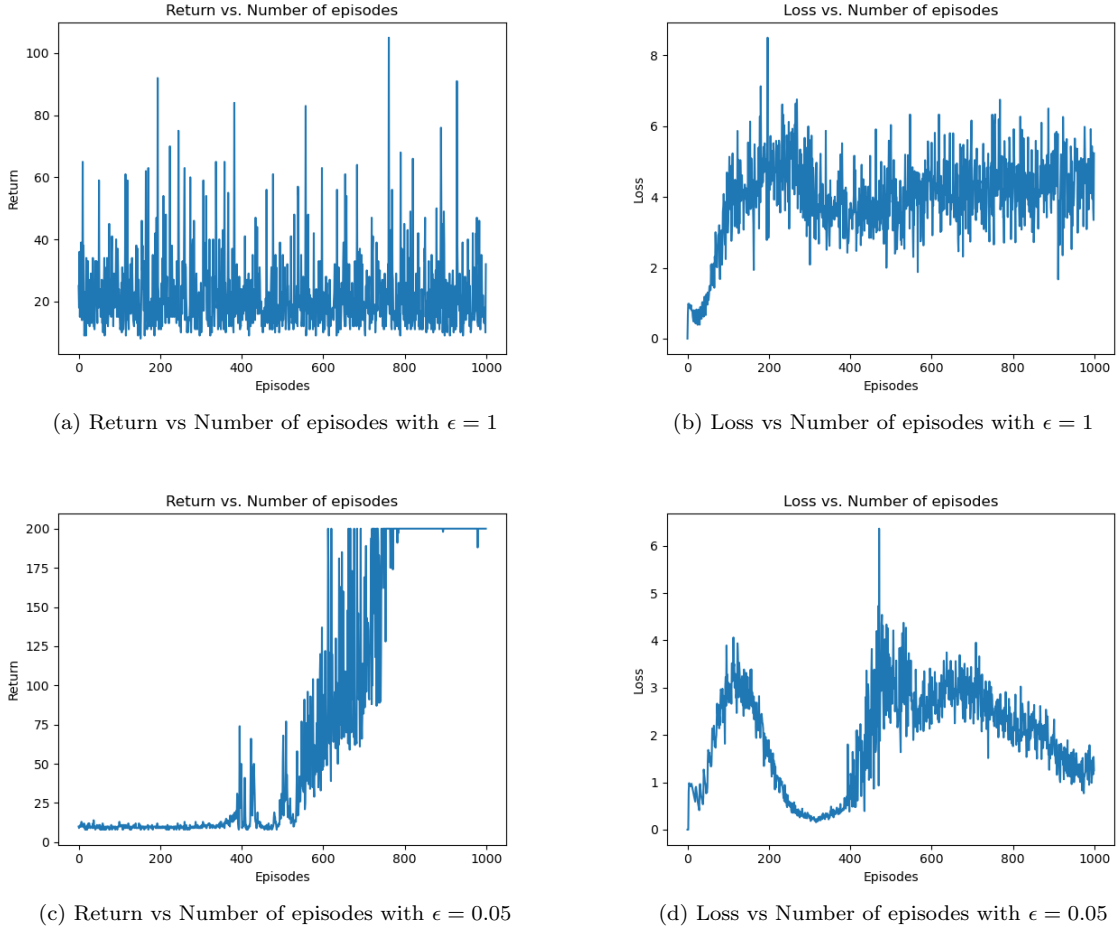


Figure 4: Plots for losses and returns vs Number of episodes for exploration and exploitation respectively

From the plots in the figures [4a, 4b], i.e., $\epsilon = 1$, the loss did not varied significantly and the returns fluctuated a lot because the action was always picked up randomly and this would account to this behaviour. On the contrary, in the figures [4c, 4d], when $\epsilon = 0.05$, after around 800 episodes the return was constant (200) and the loss was reducing after around 600 episodes.

3 DQN for Pong

3.1 Description of Pong-v0

Pong is 2D tennis table-themed arcade video game, manufactured by Atari and originally released in 1972. In Pong, one player scores if the ball passes by the other player. An episode is over when one of the players reaches 21 points. In the OpenAI Gym framework version of Pong², the agent is displayed on the right and the enemy on the left.

- The observation from the environment is an RGB-image of size (210, 160, 3) with 128 possible colors for each pixel.
- The agent has 6 possible actions (['NOOP', 'FIRE', 'RIGHT', 'LEFT', 'RIGHTFIRE', 'LEFTFIRE']) but 3 of them are redundant: actions 2/4 and 3/5 makes the paddle go down and up respectively, while action 0/1 do nothing.

²<https://gym.openai.com/envs/Pong-v0/>

- The agent receives a +1 reward if the ball went past the opponent, a -1 reward if it missed the ball, and 0 otherwise.

3.2 Implementation

The environment was Atari preprocessed, with a frame skip parameter 4 and the image was converted into grayscale image for further processing. Now each frame (observation) was converted into a pytorch tensor and proceeds into the further operations. An observation stack of size 4 was implemented where the latest observation was appended at the last index and the rest of the implementation proceeds just as in CartPole-v0 implementation.

As mentioned above some of the actions were redundant and so the actions were mapped from 0 to 2 and 1 to 3.

3.3 Results and discussion

We aimed to train the model for 10000 episodes and used the default hyperparameters to train the DQN. The hyperparameters are described in the table below.

Hyperparameter	Value	Description
Screen size	84	Resizing the frame.
Grayscale Obs	True	Returns grayscale image.
Frame Skip	1	Frequency at which agent experiences game.
Noop Max	30	Maximum number of no-ops.
Scale Obs	True	Returns the observation in range $[0, 1]$.
Memory size	10000	Maximum size of the replay buffer.
Observation stack	4	Size of observation stack.
Batch size	32	Size of mini-batches sampled from the replay buffer.
Target update frequency	1000	The frequency (measured in time steps) with which the target network is updated.
Train frequency	4	The frequency (measured in time steps) with which the network is trained.
Gamma	0.99	Discounted factor gamma used in the Q-learning update.
Learning rate	0.0001	The learning rate used by Adam optimizer.
Epsilon start	1	Initial value of ϵ -greedy exploration.
Epsilon end	0.01	Final value of ϵ -greedy exploration.
Anneal length	10^6	Value to compute the number of steps to anneal epsilon for <code>eps_step = (eps_start - eps_end) / anneal_length</code>
No. of actions	2	Number of actions the agent takes

The plots for training losses and returns during the course of training are below.

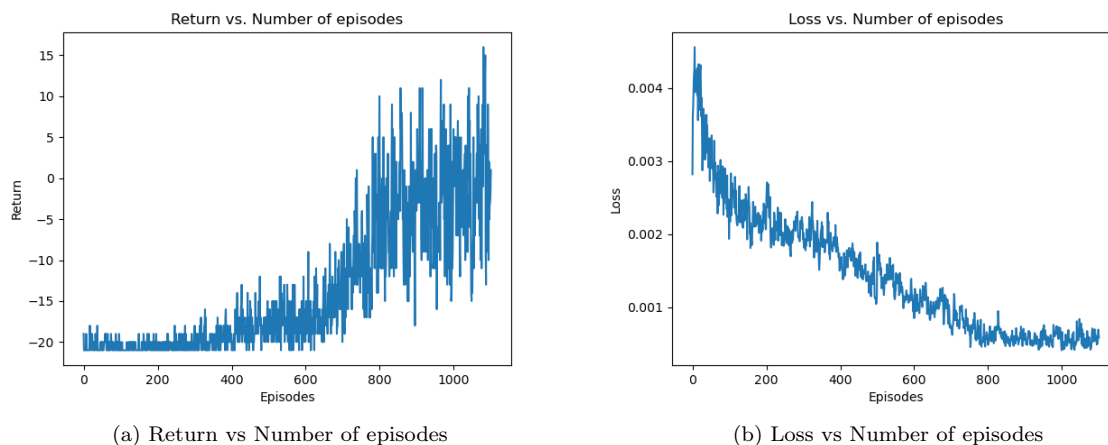


Figure 5: Losses and returns of Pong agent over 1200 episodes

At the end of 1200 episodes, our agent was able to reach a best mean return of 15.6, which is a decent value. The loss is decreasing over the course of the training, going from 0.004 to less than 0.001, suggesting that our agent is learning. Looking at the return graph, we can see that the training is becoming more and more unstable after 600 episodes: the curve fluctuates a lot. This illustrates the fact that using a neural network in reinforcement learning adds instability to the learning.

There were some few challenges we have faced while running this program:

- While testing the program, we used GPU provided in Google Colab but unfortunately GPU limit usage got exceeded forcing us to work on our local machines which took almost 8 hours to run 1000 episodes.
- We aimed to run for 10000 episodes but after certain number of episodes our systems became unresponsive and hence we ran for 1200 episodes and we attained a decent best mean return of 15.6 which had a significant improvement in return at around 300 episodes while training.

4 Appendix

```

1  """
2  In this file, you may edit the hyperparameters used for different
   environments.
3
4  memory_size: Maximum size of the replay memory.
5  n_episodes: Number of episodes to train for.
6  batch_size: Batch size used for training DQN.
7  target_update_frequency: How often to update the target network.
8  train_frequency: How often to train the DQN.
9  gamma: Discount factor.
10 lr: Learning rate used for optimizer.
11 eps_start: Starting value for epsilon (linear annealing).
12 eps_end: Final value for epsilon (linear annealing).
13 anneal_length: How many steps to anneal epsilon for.
14 n_actions: The number of actions can easily be accessed with env.
   action_space.n, but we do
15     some manual engineering to account for the fact that Pong has
   duplicate actions.
16 """

```



```

17
18 # Hyperparameters for CartPole-v0
19 CartPole = {
20     'memory_size': 50000,
21     'n_episodes': 1000,
22     'batch_size': 32,
23     'target_update_frequency': 100,
24     'train_frequency': 1,
25     'gamma': 0.95,
26     'lr': 1e-4,
27     'eps_start': 1,
28     'eps_end': 0.05,
29     'anneal_length': 10**4,
30     'n_actions': 2,
31 }
32
33 # Hyperparameters for Pong-v0
34 Pong = {
35     'obs_stack_size': 4,
36     'memory_size': 10000,
37     'n_episodes': 10000,
38     'batch_size': 32,
39     'target_update_frequency': 1000,
40     'train_frequency': 4,
41     'gamma': 0.99,
42     'lr': 1e-4,
43     'eps_start': 1.0,
44     'eps_end': 0.01,
45     'anneal_length': 10**6,
46     'n_actions': 2,
47 }

```

Listing 1: config.py

```

1 import random
2
3 import gym
4 import torch
5 import torch.nn as nn
6 import torch.nn.functional as F
7
8 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
9
10
11 class ReplayMemory:
12     def __init__(self, capacity):
13         self.capacity = capacity
14         self.memory = []
15         self.position = 0
16
17     def __len__(self):
18         return len(self.memory)
19
20     def push(self, obs, action, next_obs, reward):
21         if len(self.memory) < self.capacity:
22             self.memory.append(None)
23
24         self.memory[self.position] = (obs, action, next_obs, reward)
25         self.position = (self.position + 1) % self.capacity

```

```

26
27     def sample(self, batch_size):
28         """
29         Samples batch_size transitions from the replay memory and
30         returns a tuple
31             (obs, action, next_obs, reward)
32         """
33         sample = random.sample(self.memory, batch_size)
34         return tuple(zip(*sample)) # ((obs1, obs2), (action1, action2),
35         ...)
36
37 class DQN(nn.Module):
38     def __init__(self, env_name, env_config):
39         super(DQN, self).__init__()
40
41         self.env_name = env_name
42
43         # Save hyperparameters needed in the DQN class.
44         self.batch_size = env_config["batch_size"]
45         self.gamma = env_config["gamma"]
46         self.eps_start = env_config["eps_start"]
47         self.eps_end = env_config["eps_end"]
48         self.anneal_length = env_config["anneal_length"]
49         self.n_actions = env_config["n_actions"]
50
51         if self.env_name == 'CartPole-v0':
52             self.fc1 = nn.Linear(4, 256)
53             self.fc2 = nn.Linear(256, self.n_actions)
54         elif self.env_name == 'Pong-v0':
55             self.conv1 = nn.Conv2d(4, 32, kernel_size=8, stride=4,
padding=0)
56             self.conv2 = nn.Conv2d(32, 64, kernel_size=4, stride=2,
padding=0)
57             self.conv3 = nn.Conv2d(64, 64, kernel_size=3, stride=1,
padding=0)
58             self.fc1 = nn.Linear(3136, 512)
59             self.fc2 = nn.Linear(512, self.n_actions)
60
61         self.relu = nn.ReLU()
62         self.flatten = nn.Flatten()
63
64     def forward(self, x):
65         """Runs the forward pass of the NN depending on architecture."""
66
67         if self.env_name == 'CartPole-v0':
68             x = self.relu(self.fc1(x))
69             x = self.fc2(x)
70         elif self.env_name == 'Pong-v0':
71             x = self.relu(self.conv1(x))
72             x = self.relu(self.conv2(x))
73             x = self.relu(self.conv3(x))
74             x = self.flatten(x)
75             x = self.relu(self.fc1(x))
76             x = self.fc2(x)
77
78         return x

```

```

79     def act(self, observation, epsilon, exploit=False):
80         """Selects an action with an epsilon-greedy exploration strategy
81         .
82         Args:
83             observation (tensor): observation tensor of shape [
batch_size, state_dimension] (e.g. [32, 4])
84             epsilon (int)
85             exploit (bool): if True, acts greedily
86
87         Returns:
88             actions (tensor): tensor of actions according to the epsilon
-greedy strategy, of shape [batch_size, 1]
89         """
90
91         batch_size = observation.shape[0]
92         rand_value = random.random()
93
94         if exploit or rand_value > epsilon:
95             # Choose the action which gives the largest Q-values for
each observation
96             with torch.no_grad():
97                 output = self.forward(observation)
98                 actions = torch.argmax(output, dim=1)
99
100                 if self.env_name == "Pong-v0":
101                     if actions.item() == 0:
102                         actions = torch.tensor([2]) # right
103                     else:
104                         actions = torch.tensor([3]) # left
105                 else:
106                     # Choose a random action for each observation
107                     if self.env_name == "CartPole-v0":
108                         random_actions = [random.randrange(self.n_actions) for _
in range(batch_size)]
109                     if self.env_name == "Pong-v0":
110                         random_actions = [random.randrange(2, 4) for _ in range(
batch_size)]
111                     actions = torch.tensor(random_actions)
112                 return actions.unsqueeze(1)
113
114     def optimize(dqn, target_dqn, memory, optimizer):
115         """This function samples a batch from the replay buffer and
optimizes the Q-network."""
116
117         # If we don't have enough transitions stored yet, we don't train.
118         if len(memory) < dqn.batch_size:
119             return 0
120
121         # Sample a batch from the replay memory and concatenate so that
there are
122         # four tensors in total: observations, actions, next observations
and rewards.
123         observations, actions, next_observations, rewards = memory.sample(
dqn.batch_size)
124
125         # Transform every tuple into tensors.
126         observations = torch.cat(observations)

```

```

127     non_terminal_next_observations = [next_obs for next_obs in
128     next_observations if next_obs is not None]
129     non_terminal_next_observations = torch.cat(
130     non_terminal_next_observations).float()
131
132     actions = torch.stack(actions, dim=0)
133     actions = torch.unsqueeze(actions, 1)
134
135     rewards = torch.stack(rewards, dim=0)
136     rewards = torch.unsqueeze(rewards, 1)
137
138     # Compute the current estimates of the Q-values for each state-
139     # action pair (s,a).
140     output = dqn.forward(observations)
141     output = output.to(device)
142     if dqn.env_name == "CartPole-v0":
143         q_values = torch.gather(input=output, index=actions, dim=1)
144     if dqn.env_name == "Pong-v0":
145         q_values = torch.gather(input=output, index=actions - 2, dim=1)
146
147     # Compute the Q-value targets.
148     non_terminal_mask = torch.BoolTensor(list(map(lambda obs: obs is not
149     None, next_observations))) # indices for non terminal transitions
150     terminal_mask = ~non_terminal_mask # indices for terminal
151     transitions
152     target_dqn_qpred = target_dqn.forward(non_terminal_next_observations
153     )
154     target_dqn_qpred_max = torch.max(target_dqn_qpred, axis=1)[0].
155     unsqueeze(1)
156
157     q_value_targets = torch.zeros(dqn.batch_size, 1)
158     q_value_targets = q_value_targets.to(device)
159     q_value_targets[non_terminal_mask] = rewards[non_terminal_mask] +
160     dqn.gamma * target_dqn_qpred_max
161     q_value_targets[terminal_mask] = rewards[terminal_mask]
162
163     # Compute loss.
164     loss = F.mse_loss(q_values, q_value_targets)
165
166     # Perform gradient descent.
167     optimizer.zero_grad()
168
169     loss.backward()
170     optimizer.step()
171
172     return loss.item()

```

Listing 2: dqn.py

```

1 import argparse
2 import random
3
4 import gym
5 import torch
6 import torch.nn as nn
7
8 import config
9 from utils import preprocess

```

```

10 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
11
12
13 parser = argparse.ArgumentParser()
14 parser.add_argument('--env', default='Pong-v0', choices=['CartPole-v0',
15 'Pong-v0'])
16 parser.add_argument('--path', type=str, help='Path to stored DQN model.'
17 )
18 parser.add_argument('--n_eval_episodes', type=int, default=1, help='
19 Number of evaluation episodes.', nargs='?')
20 parser.add_argument('--render', dest='render', action='store_true', help
21 ='Render the environment.')
22 parser.add_argument('--save_video', dest='save_video', action='
23 store_true', help='Save the episodes as video.')
24 parser.set_defaults(render=False)
25 parser.set_defaults(save_video=False)
26
27 # Hyperparameter configurations for different environments. See config.
28 py.
29 ENV_CONFIGS = {
30     'CartPole-v0': config.CartPole,
31     'Pong-v0': config.Pong
32 }
33
34 def evaluate_policy(dqn, env, env_config, args, eps_start, n_episodes,
35 discounted=False, render=False, verbose=False):
36     """Runs {n_episodes} episodes to evaluate current policy."""
37     total_return = 0
38
39     for i in range(n_episodes):
40         obs = preprocess(env.reset(), env=args.env).unsqueeze(0)
41         if args.env == "Pong-v0":
42             obs_stack = torch.cat(env_config["obs_stack_size"] * [obs]).
43             unsqueeze(0).to(device)
44
45         done = False
46         rewards_list = []
47
48         while not done:
49             if render:
50                 env.render()
51
52             if args.env == "Pong-v0":
53                 action = dqn.act(obs_stack, eps_start, exploit=True).
54                 item()
55             else:
56                 action = dqn.act(obs, eps_start, exploit=True).item()
57
58             obs, reward, done, info = env.step(action)
59             obs = preprocess(obs, env=args.env).unsqueeze(0)
60             if args.env == "Pong-v0":
61                 obs_stack = torch.cat((obs_stack[:, 1:, ...], obs.
62                 unsqueeze(1)), dim=1).to(device)
63                 obs_stack = preprocess(obs_stack, env=args.env)
64
65             rewards_list.append(reward)
66
67         if discounted:

```

```

58         episode_return = sum([env_config['gamma']**t * rewards_list[
    t] for t in range(len(rewards_list))])
59     else:
60         episode_return = sum(rewards_list)
61     total_return += episode_return
62
63     if verbose:
64         print(f'Finished episode {i+1} with a total return of {
    episode_return}')
65
66
67     return total_return / n_episodes
68
69 if __name__ == '__main__':
70     args = parser.parse_args()
71
72     # Initialize environment and config
73     env = gym.make(args.env)
74     env_config = ENV_CONFIGS[args.env]
75
76     if args.save_video:
77         env = gym.wrappers.Monitor(env, './video/', video_callable=
    lambda episode_id: True, force=True)
78
79     # Load model from provided path.
80     dqn = torch.load(args.path, map_location=torch.device('cpu'))
81     dqn.eval()
82
83     mean_return = evaluate_policy(dqn, env, env_config, args, dqn.
    eps_start, args.n_eval_episodes, render=args.render and not args.
    save_video, verbose=True)
84     print(f'The policy got a mean return of {mean_return} over {args.
    n_eval_episodes} episodes.')
85
86     env.close()

```

Listing 3: evaluate.py

```

1  import argparse
2
3  import gym
4  from gym.wrappers import AtariPreprocessing
5  import torch
6  import torch.nn as nn
7
8  import config
9  from utils import preprocess, init_weights, plot_result
10 from evaluate import evaluate_policy
11 from dqn import DQN, ReplayMemory, optimize
12
13 import matplotlib.pyplot as plt
14 import pickle
15
16 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
17
18 parser = argparse.ArgumentParser()
19 parser.add_argument('--env', default='CartPole-v0', choices=['CartPole-
    v0', 'Pong-v0'])
20 parser.add_argument('--evaluate_freq', type=int, default=25, help='How

```

```

    often to run evaluation.', nargs='?')
21 parser.add_argument('--evaluation_episodes', type=int, default=5, help='
    Number of evaluation episodes.', nargs='?')
22 parser.add_argument('--model_name', default='Pong-v0', help='Model to
    load if present.', nargs='?')
23
24 # Hyperparameter configurations for different environments. See config.
    py.
25 ENV_CONFIGS = {
26     'CartPole-v0': config.CartPole,
27     'Pong-v0': config.Pong
28 }
29
30 if __name__ == '__main__':
31     args = parser.parse_args()
32
33     # Initialize environment and config.
34     env = gym.make(args.env)
35     if args.env == "Pong-v0":
36         env = AtariPreprocessing(env, screen_size=84, grayscale_obs=True
    , frame_skip=1, noop_max=30, scale_obs=True)
37     env_config = ENV_CONFIGS[args.env]
38     params = f"lr_{env_config['lr']}_gamma_{env_config['gamma']}_tg_{
    env_config['target_update_frequency']}_ann_{env_config['anneal_length
    ']}]"
39
40     # Load Model.
41     try:
42         model_name = f"{args.env}_best_{params}.pt"
43         target_model_name = f"{args.env}_target_{params}.pt"
44         dqn = torch.load(f'models/{model_name}')
45         target_dqn = torch.load(f'models/{target_model_name}')
46         print("Loading model...")
47     except:
48         print("Initialize model...")
49
50     # Initialize deep Q-networks.
51     dqn = DQN(env_name=args.env, env_config=env_config).to(device)
52     dqn.apply(init_weights)
53
54     # Create and initialize target Q-network.
55     target_dqn = DQN(env_name=args.env, env_config=env_config).to(
    device)
56     target_dqn.load_state_dict(dqn.state_dict())
57
58     # Load results pickle file to add new results from the pretrained
    model.
59     results_file = f"results/{args.env}/{args.env}_results_{params}.pkl"
60     try:
61         with open(results_file, 'rb') as file:
62             results_dict = pickle.load(file)
63             print("Load pickle file")
64             return_list = results_dict['return']
65             eval_return_list = results_dict['eval_return']
66             loss_list = results_dict['loss']
67             step_list = results_dict['step']
68             epsilon_list = results_dict['epsilon']
69             dqn.eps_start = results_dict['eps_start']

```

```

70     best_mean_return = results_dict['best_mean_return']
71     memory = results_dict['memory']
72     print(f"best_mean_return = {best_mean_return}")
73 except:
74     # Initialize lists to keep track of episodes' loss, return,
75     # number of steps, and epsilon values throughout the training.
76     loss_list, return_list, eval_return_list, step_list,
77     epsilon_list = [], [], [], [], []
78
79     # Keep track of best evaluation mean return achieved so far.
80     best_mean_return = -float("Inf")
81
82     # Create replay memory.
83     memory = ReplayMemory(env_config['memory_size'])
84
85     # Initialize optimizer used for training the DQN. We use Adam rather
86     # than RMSProp.
87     optimizer = torch.optim.Adam(dqn.parameters(), lr=env_config['lr'])
88
89     # Compute the value for epsilon linear annealing.
90     eps = dqn.eps_start
91     eps_end = dqn.eps_end
92     anneal_length = dqn.anneal_length
93     eps_step = (eps - eps_end) / anneal_length
94
95     # Keep track of the total number of steps.
96     total_steps = 0
97
98     for episode in range(env_config['n_episodes']):
99         done = False
100
101         episode_loss = 0
102         episode_rewards = []
103         episode_steps = 0
104
105         obs = preprocess(env.reset(), env=args.env).unsqueeze(0).to(
106             device)
107         if args.env == "Pong-v0":
108             obs_stack = torch.cat(env_config["obs_stack_size"] * [obs]).
109             unsqueeze(0).to(device)
110
111         while not done:
112             total_steps += 1
113             episode_steps += 1
114
115             # Get action from DQN.
116             if args.env == "Pong-v0":
117                 action = dqn.act(obs_stack, eps).item()
118             else:
119                 action = dqn.act(obs, eps).item()
120
121             # Act in the true environment.
122             next_obs, reward, done, info = env.step(action)
123             episode_rewards.append(reward)
124
125             # Preprocess incoming observation.
126             if not done:
127                 # Preprocess the non-terminal state.

```



```

123         next_obs = preprocess(next_obs, env=args.env).unsqueeze
124         (0).to(device)
125         if args.env == "Pong-v0":
126             next_obs_stack = torch.cat((obs_stack[:, 1:, ...],
127 next_obs.unsqueeze(1)), dim=1).to(device)
128             next_obs_stack = preprocess(next_obs_stack, env=args
129 .env).to(device)
130         else:
131             # Set to None the terminal state.
132             next_obs = None
133             if args.env == "Pong-v0":
134                 next_obs_stack = None
135
136         # Add the transition to the replay memory. Everything has
137         been move to GPU if available.
138         action = torch.as_tensor(action).to(device)
139         reward = torch.as_tensor(reward).to(device)
140         if args.env == "CartPole-v0":
141             memory.push(obs, action, next_obs, reward)
142         if args.env == "Pong-v0":
143             memory.push(obs_stack, action, next_obs_stack, reward)
144
145         # Run DQN.optimize() every env_config["train_frequency"]
146         steps.
147         if total_steps % env_config["train_frequency"] == 0:
148             loss = optimize(dqn, target_dqn, memory, optimizer)
149             episode_loss += loss
150
151         # Update the target network every env_config["
152         target_update_frequency"] steps.
153         if total_steps % env_config["target_update_frequency"] == 0:
154             target_dqn.load_state_dict(dqn.state_dict())
155
156         # Update the current observation.
157         obs = next_obs
158         if args.env == "Pong-v0":
159             obs_stack = next_obs_stack
160
161         # Update epsilon after each step.
162         epsilon_list.append(eps)
163         if eps > eps_end:
164             eps -= eps_step
165
166         # Compute the episode return.
167         episode_undiscounted_return = sum(episode_rewards)
168
169         # Add results to lists.
170         loss_list.append(episode_loss/episode_steps)
171         return_list.append(episode_undiscounted_return)
172         step_list.append(episode_steps)
173
174         # Evaluate the current agent.
175         if episode % args.evaluate_freq == 0:
176             mean_return = evaluate_policy(dqn, env, env_config, args,
177 eps, n_episodes=args.evaluation_episodes)
178             eval_return_list.append(mean_return)
179
180         print(f'Episode {episode}/{env_config["n_episodes"]}: {

```

```

mean_return}')
174
175     # Save current agent if it has the best performance so far.
176     if mean_return >= best_mean_return:
177         best_mean_return = mean_return
178
179         print('Best performance so far! Saving model.')
180         torch.save(dqn, f'models/{args.env}_best_{params}.pt')
181         torch.save(target_dqn, f'models/{args.env}_target_{
params}.pt')
182
183     # Save results.
184     results_dict = {
185         'return': return_list,
186         'eval_return': eval_return_list,
187         'loss': loss_list,
188         'step': step_list,
189         'epsilon': epsilon_list,
190         'eps_start': eps,
191         'best_mean_return': best_mean_return,
192         'memory': memory
193     }
194     print("Save pickle file")
195     with open(results_file, 'wb') as file:
196         pickle.dump(results_dict, file)
197
198     # Plot
199     # x_axis = range(len(results_dict['return']))
200     # plot_result(x_axis, results_dict['return'], "Episodes", "
Return", title="Return vs. Number of episodes", save_path=f"results/{
args.env}", env_name=args.env, params=params)
201     # plot_result(range(len(results_dict['eval_return'])),
results_dict['eval_return'], "Episodes", "Return", title="Evaluation
Return every 25 episodes", save_path=f"results/{args.env}", env_name=
args.env, params=params)
202     # plot_result(x_axis, results_dict['loss'], "Episodes", "
Loss", title="Loss vs. Number of episodes", save_path=f"results/{args
.env}", env_name=args.env, params=params)
203     # plot_result(x_axis, results_dict['step'], "Episodes", "
Number of steps", title="Number of steps per episode", save_path=f"
results/{args.env}", env_name=args.env, params=params)
204     # plot_result(range(len(results_dict['epsilon'])),
results_dict['epsilon'], "Steps", "Epsilon", title="Epsilon values at
each step of the training", save_path=f"results/{args.env}",
env_name=args.env, params=params)
205
206     # Close environment after training is completed.
207     env.close()
208
209     # Save results.
210     results_dict = {
211         'return': return_list,
212         'eval_return': eval_return_list,
213         'loss': loss_list,
214         'step': step_list,
215         'epsilon': epsilon_list,
216         'eps_start': eps,
217         'best_mean_return': best_mean_return,

```

```

218         'memory': memory
219     }
220     print("Save pickle file")
221     with open(results_file, 'wb') as file:
222         pickle.dump(results_dict, file)
223
224     # Plot.
225     x_axis = range(len(results_dict['return']))
226     plot_result(x_axis, results_dict['return'], "Episodes", "Return",
227               title="Return vs. Number of episodes", save_path=f"results/{args.env}",
228               env_name=args.env, params=params)
229     plot_result(range(len(results_dict['eval_return'])), results_dict['eval_return'], "Episodes", "Return", title="Evaluation Return every 25 episodes", save_path=f"results/{args.env}", env_name=args.env, params=params)
230     plot_result(x_axis, results_dict['loss'], "Episodes", "Loss", title="Loss vs. Number of episodes", save_path=f"results/{args.env}", env_name=args.env, params=params)
231     plot_result(x_axis, results_dict['step'], "Episodes", "Number of steps", title="Number of steps per episode", save_path=f"results/{args.env}", env_name=args.env, params=params)
232     plot_result(range(len(results_dict['epsilon'])), results_dict['epsilon'], "Steps", "Epsilon", title="Epsilon values at each step of the training", save_path=f"results/{args.env}", env_name=args.env, params=params)

```

Listing 4: train.py

```

1  import random
2
3  import gym
4  from gym.wrappers import AtariPreprocessing
5  import torch
6  import torch.nn as nn
7
8  import matplotlib.pyplot as plt
9
10 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
11
12 def preprocess(obs, env):
13     """Performs necessary observation preprocessing."""
14     if env in ['CartPole-v0', 'Pong-v0']:
15         return torch.tensor(obs, device=device).float()
16     else:
17         raise ValueError('Please add necessary observation preprocessing instructions to preprocess() in utils.py.')
18
19 def init_weights(model):
20     if isinstance(model, torch.nn.Linear):
21         torch.nn.init.normal_(model.weight, mean=0, std=0.01)
22         torch.nn.init.zeros_(model.bias)
23
24 def plot_result(x, y, xlabel, ylabel, title, save_path, env_name, params):
25     plt.plot(x, y)
26     plt.title(title)
27     plt.xlabel(xlabel)
28     plt.ylabel(ylabel)
29     plt.savefig(f"{save_path}/{env_name}_{title}_{params}.png")

```

```
plt.show()
```

Listing 5: utils.py

References

- [1] Volodymyr Mnih et al. “Human-level control through deep reinforcement learning”. In: *Nature* 518.7540 (Feb. 2015), pp. 529–533. ISSN: 00280836. URL: <http://dx.doi.org/10.1038/nature14236>.
- [2] Volodymyr Mnih et al. “Playing Atari with Deep Reinforcement Learning”. In: *CoRR* abs/1312.5602 (2013). arXiv: [1312.5602](https://arxiv.org/abs/1312.5602). URL: <http://arxiv.org/abs/1312.5602>.