

#Tutorial 5 - DQN

Please follow this tutorial to understand the structure (code) of DQN algorithm.

References:

Please follow [Human-level control through deep reinforcement learning \(https://www.nature.com/articles/nature14236\)](https://www.nature.com/articles/nature14236) for the original publication as well as the psuedocode. Watch Prof. Ravi's lectures on moodle or nptel for further understanding of the core concepts. Contact the TAs for further resources if needed.

In [2]:

```
1 '''
2 Installing packages for rendering the game on Colab
3 '''
4
5 !pip install gym pyvirtualdisplay > /dev/null 2>&1
6 !apt-get install -y xvfb python-opengl ffmpeg > /dev/null 2>&1
7 !apt-get update > /dev/null 2>&1
8 !apt-get install cmake > /dev/null 2>&1
9 !pip install --upgrade setuptools 2>&1
10 !pip install ez_setup > /dev/null 2>&1
11 !pip install gym[atari] > /dev/null 2>&1
12 !pip install git+https://github.com/tensorflow/docs > /dev/null 2>&1
13 !pip install gym[classic_control]
```

Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (76.0.0)
Requirement already satisfied: gym[classic_control] in /usr/local/lib/python3.11/dist-packages (0.26.2)
Requirement already satisfied: numpy>=1.18.0 in /usr/local/lib/python3.11/dist-packages (from gym[classic_control]) (1.26.4)
Requirement already satisfied: cloudpickle>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from gym[classic_control]) (3.1.1)
Requirement already satisfied: gym_notices>=0.0.4 in /usr/local/lib/python3.11/dist-packages (from gym[classic_control]) (0.0.8)
Collecting pygame==2.1.0 (from gym[classic_control])
Using cached pygame-2.1.0.tar.gz (5.8 MB)
error: subprocess-exited-with-error

× python setup.py egg_info did not run successfully.
 exit code: 1
 ↳ See above for output.

note: This error originates from a subprocess, and is likely not a problem with pip.
Preparing metadata (setup.py) ... error
error: metadata-generation-failed

× Encountered error while generating package metadata.
 ↳ See above for output.

note: This is an issue with the package mentioned above, not pip.
hint: See above for details.

In [3]:

```
1 '''
2 A bunch of imports, you don't have to worry about these
3 '''
4
5 import numpy as np
6 import random
7 import torch
8 import torch.nn as nn
9 import torch.nn.functional as F
10 from collections import namedtuple, deque
11 import torch.optim as optim
12 import datetime
13 import gymnasium as gym
14 from gym.wrappers.record_video import RecordVideo
15 import glob
16 import io
17 import base64
18 import matplotlib.pyplot as plt
19 from IPython.display import HTML
20 from pyvirtualdisplay import Display
21 import tensorflow as tf
22 from IPython import display as ipythondisplay
23 from PIL import Image
24 import tensorflow_probability as tfp
```

```
In [4]: 1 '''
2 Please refer to the first tutorial for more details on the specifics of environments
3 We've only added important commands you might find useful for experiments.
4 '''
5
6 '''
7 List of example environments
8 (Source - https://gym.openai.com/envs/#classic_control)
9
10 'Acrobot-v1'
11 'Cartpole-v1'
12 'MountainCar-v0'
13 '''
14 env = gym.make('CartPole-v1')
15
16 state_shape = env.observation_space.shape[0]
17 no_of_actions = env.action_space.n
18
19 print(state_shape)
20 print(no_of_actions)
21 print(env.action_space.sample())
22 print("----")
23
24 '''
25 # Understanding State, Action, Reward Dynamics
26
27 The agent decides an action to take depending on the state.
28
29 The Environment keeps a variable specifically for the current state.
30 - Everytime an action is passed to the environment, it calculates the new state and updates the current state var
31 - It returns the new current state and reward for the agent to take the next action
32
33 '''
34
35 state = env.reset()
36 ''' This returns the initial state (when environment is reset) '''
37
38 print(state)
39 print("----")
40
41 action = env.action_space.sample()
42 ''' We take a random action now '''
43
44 print(action)
45 print("----")
46
47 next_state, reward, done, _ , info = env.step(action)
48 ''' env.step is used to calculate new state and obtain reward based on old state and action taken '''
49
50 print(next_state)
51 print(reward)
52 print(done)
53 print(info)
54 print("----")
55
```

```
4
2
0
----
(array([-0.0410322 , -0.04988823,  0.01525484,  0.01642654], dtype=float32), {})
```

```
----
0
----
[-0.04202996 -0.2452256   0.01558337  0.31388324]
1.0
False
{}
----
```

DQN

Using NNs as substitutes isn't something new. It has been tried earlier, but the 'human control' paper really popularised using NNs by providing a few stability ideas (Q-Targets, Experience Replay & Truncation). The 'Deep-Q Network' (DQN) Algorithm can be broken down into having the following components.

Q-Network:

The neural network used as a function approximator is defined below

In [5]:

```
1  '''
2  ### Q Network & Some 'hyperparameters'
3
4  QNetwork1:
5  Input Layer - 4 nodes (State Shape) \
6  Hidden Layer 1 - 128 nodes \
7  Hidden Layer 2 - 64 nodes \
8  Output Layer - 2 nodes (Action Space) \
9  Optimizer - zero_grad()
10 '''
11
12 import torch
13 import torch.nn as nn
14 import torch.nn.functional as F
15
16
17 '''
18 Bunch of Hyper parameters (Which you might have to tune later)
19 '''
20 BUFFER_SIZE = int(1e5) # replay buffer size
21 BATCH_SIZE = 64 # minibatch size
22 GAMMA = 0.99 # discount factor
23 LR = 5e-4 # Learning rate
24 UPDATE_EVERY = 20 # how often to update the network (When Q target is present)
25
26
27 class QNetwork1(nn.Module):
28
29     def __init__(self, state_size, action_size, seed, fc1_units=128, fc2_units=64):
30         """Initialize parameters and build model.
31         Params
32         =====
33             state_size (int): Dimension of each state
34             action_size (int): Dimension of each action
35             seed (int): Random seed
36             fc1_units (int): Number of nodes in first hidden layer
37             fc2_units (int): Number of nodes in second hidden layer
38         """
39         super(QNetwork1, self).__init__()
40         self.seed = torch.manual_seed(seed)
41         self.fc1 = nn.Linear(state_size, fc1_units)
42         self.fc2 = nn.Linear(fc1_units, fc2_units)
43         self.fc3 = nn.Linear(fc2_units, action_size)
44
45     def forward(self, state):
46         """Build a network that maps state -> action values."""
47         x = F.relu(self.fc1(state))
48         x = F.relu(self.fc2(x))
49         return self.fc3(x)
```

Replay Buffer:

Recall why we use such a technique.

In [6]:

```
1 import random
2 import torch
3 import numpy as np
4 from collections import deque, namedtuple
5
6 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
7
8 class ReplayBuffer:
9     """Fixed-size buffer to store experience tuples."""
10
11     def __init__(self, action_size, buffer_size, batch_size, seed):
12         """Initialize a ReplayBuffer object.
13
14         Params
15         =====
16             action_size (int): dimension of each action
17             buffer_size (int): maximum size of buffer
18             batch_size (int): size of each training batch
19             seed (int): random seed
20         """
21         self.action_size = action_size
22         self.memory = deque(maxlen=buffer_size)
23         self.batch_size = batch_size
24         self.experience = namedtuple("Experience", field_names=["state", "action", "reward", "next_state", "done"])
25         self.seed = random.seed(seed)
26
27     def add(self, state, action, reward, next_state, done):
28         """Add a new experience to memory."""
29         e = self.experience(state, action, reward, next_state, done)
30         self.memory.append(e)
31
32     def sample(self):
33         """Randomly sample a batch of experiences from memory."""
34         experiences = random.sample(self.memory, k=self.batch_size)
35
36         states = torch.from_numpy(np.vstack([e.state for e in experiences if e is not None])).float().to(device)
37         actions = torch.from_numpy(np.vstack([e.action for e in experiences if e is not None])).long().to(device)
38         rewards = torch.from_numpy(np.vstack([e.reward for e in experiences if e is not None])).float().to(device)
39         next_states = torch.from_numpy(np.vstack([e.next_state for e in experiences if e is not None])).float().to(device)
40         dones = torch.from_numpy(np.vstack([e.done for e in experiences if e is not None]).astype(np.uint8)).float().to(device)
41
42         return (states, actions, rewards, next_states, dones)
43
44     def __len__(self):
45         """Return the current size of internal memory."""
46         return len(self.memory)
```

Tutorial Agent Code:

```
In [7]: 1 class TutorialAgent():
2
3     def __init__(self, state_size, action_size, seed):
4
5         ''' Agent Environment Interaction '''
6         self.state_size = state_size
7         self.action_size = action_size
8         self.seed = random.seed(seed)
9
10        ''' Q-Network '''
11        self.qnetwork_local = QNetwork1(state_size, action_size, seed).to(device)
12        self.qnetwork_target = QNetwork1(state_size, action_size, seed).to(device)
13        self.optimizer = optim.Adam(self.qnetwork_local.parameters(), lr=LR)
14
15        ''' Replay memory '''
16        self.memory = ReplayBuffer(action_size, BUFFER_SIZE, BATCH_SIZE, seed)
17
18        ''' Initialize time step (for updating every UPDATE_EVERY steps)          -Needed for Q Targets '''
19        self.t_step = 0
20
21    def step(self, state, action, reward, next_state, done):
22
23        ''' Save experience in replay memory '''
24        self.memory.add(state, action, reward, next_state, done)
25
26        ''' If enough samples are available in memory, get random subset and learn '''
27        if len(self.memory) >= BATCH_SIZE:
28            experiences = self.memory.sample()
29            self.learn(experiences, GAMMA)
30
31        """ +Q TARGETS PRESENT """
32        ''' Updating the Network every 'UPDATE_EVERY' steps taken '''
33        self.t_step = (self.t_step + 1) % UPDATE_EVERY
34        if self.t_step == 0:
35
36            self.qnetwork_target.load_state_dict(self.qnetwork_local.state_dict())
37
38    def act(self, state, eps=0.):
39
40        state = torch.from_numpy(state).float().unsqueeze(0).to(device)
41        self.qnetwork_local.eval()
42        with torch.no_grad():
43            action_values = self.qnetwork_local(state)
44        self.qnetwork_local.train()
45
46        ''' Epsilon-greedy action selection (Already Present) '''
47        if random.random() > eps:
48            return np.argmax(action_values.cpu().data.numpy())
49        else:
50            return random.choice(np.arange(self.action_size))
51
52    def learn(self, experiences, gamma):
53        """ +E EXPERIENCE REPLAY PRESENT """
54        states, actions, rewards, next_states, dones = experiences
55
56        ''' Get max predicted Q values (for next states) from target model'''
57        Q_targets_next = self.qnetwork_target(next_states).detach().max(1)[0].unsqueeze(1)
58
59        ''' Compute Q targets for current states '''
60        Q_targets = rewards + (gamma * Q_targets_next * (1 - dones))
61
62        ''' Get expected Q values from local model '''
63        Q_expected = self.qnetwork_local(states).gather(1, actions)
64
65        ''' Compute loss '''
66        loss = F.mse_loss(Q_expected, Q_targets)
67
68        ''' Minimize the loss '''
69        self.optimizer.zero_grad()
70        loss.backward()
71
72        ''' Gradient Clipping '''
73        """ +T TRUNCATION PRESENT """
74        for param in self.qnetwork_local.parameters():
75            param.grad.data.clamp_(-1, 1)
76
77        self.optimizer.step()
```

Here, we present the DQN algorithm code.

```
In [7]: 1 ''' Defining DQN Algorithm '''
2
3 state_shape = env.observation_space.shape[0]
4 action_shape = env.action_space.n
5
6
7 def dqn(n_episodes=10000, max_t=1000, eps_start=1.0, eps_end=0.01, eps_decay=0.995):
8
9     scores_window = deque(maxlen=100)
10     ''' last 100 scores for checking if the avg is more than 195 '''
11
12     eps = eps_start
13     ''' initialize epsilon '''
14
15     for i_episode in range(1, n_episodes+1):
16         state, _ = env.reset()
17         score = 0
18         for t in range(max_t):
19             action = agent.act(state, eps)
20             next_state, reward, done, _, _ = env.step(action)
21             agent.step(state, action, reward, next_state, done)
22             state = next_state
23             score += reward
24             if done:
25                 break
26
27         scores_window.append(score)
28
29         eps = max(eps_end, eps_decay*eps)
30         ''' decrease epsilon '''
31
32         print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)), end='')
33
34         if i_episode % 100 == 0:
35             print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)))
36         if np.mean(scores_window)>=195.0:
37             print('\nEnvironment solved in {:d} episodes! \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores
38                 break
39     return True
40
41 ''' Trial run to check if algorithm runs and saves the data '''
42
43 begin_time = datetime.datetime.now()
44
45 agent = TutorialAgent(state_size=state_shape,action_size = action_shape,seed = 0)
46 dqn()
47
48 time_taken = datetime.datetime.now() - begin_time
49
50 print(time_taken)
```

```
Episode 100      Average Score: 36.88
Episode 200      Average Score: 160.92
Episode 300      Average Score: 60.54
Episode 400      Average Score: 11.77
Episode 500      Average Score: 40.99
Episode 600      Average Score: 13.68
Episode 700      Average Score: 117.94
Episode 800      Average Score: 25.25
Episode 900      Average Score: 147.36
Episode 1000     Average Score: 167.90
Episode 1100     Average Score: 164.52
Episode 1200     Average Score: 162.53
Episode 1300     Average Score: 174.43
Episode 1354     Average Score: 197.40
Environment solved in 1354 episodes!      Average Score: 197.40
0:07:48.237709
```

Task 1a

Understand the core of the algorithm, follow the flow of data. Identify the exploration strategy used.

Task 1b

Out of the two exploration strategies discussed in class (ϵ -greedy & Softmax). Implement the strategy that's not used here.

Task 1c

How fast does the agent 'solve' the environment in terms of the number of episodes? (Cartpole-v1 defines "solving" as getting average reward of 195.0 over 100 consecutive trials)

How 'well' does the agent learn? (reward plot?) The above two are some 'evaluation metrics' you can use to comment on the performance of an algorithm.

Please compare DQN (using ϵ -greedy) with DQN (using softmax). Think along the lines of 'no. of episodes', 'reward plots', 'compute time', etc. and add a few comments.

Submission Steps

Task 1: Add a text cell with the answer.

Task 2: Add a code cell below task 1 solution and use 'Tutorial Agent Code' to build your new agent (with a different exploration strategy).

Task 3: Add a code cell below task 2 solution running both the agents to solve the CartPole v-1 environment and add a new text cell below it with your inferences.

Task 1 Solution

The Exploration strategy used in the given TutorialAgent class is Epsilon greedy.

Task 2 Solution

In [8]:

```

1  class TutorialAgent_SoftMax():
2
3      def __init__(self, state_size, action_size, seed):
4
5          ''' Agent Environment Interaction '''
6          self.state_size = state_size
7          self.action_size = action_size
8          self.seed = random.seed(seed)
9
10         ''' Q-Network '''
11         self.qnetwork_local = QNetwork1(state_size, action_size, seed).to(device)
12         self.qnetwork_target = QNetwork1(state_size, action_size, seed).to(device)
13         self.optimizer = optim.Adam(self.qnetwork_local.parameters(), lr=LR)
14
15         ''' Replay memory '''
16         self.memory = ReplayBuffer(action_size, BUFFER_SIZE, BATCH_SIZE, seed)
17
18         ''' Initialize time step (for updating every UPDATE_EVERY steps)          -Needed for Q Targets '''
19         self.t_step = 0
20
21     def step(self, state, action, reward, next_state, done):
22
23         ''' Save experience in replay memory '''
24         self.memory.add(state, action, reward, next_state, done)
25
26         ''' If enough samples are available in memory, get random subset and learn '''
27         if len(self.memory) >= BATCH_SIZE:
28             experiences = self.memory.sample()
29             self.learn(experiences, GAMMA)
30
31         """ +Q TARGETS PRESENT """
32         ''' Updating the Network every 'UPDATE_EVERY' steps taken '''
33         self.t_step = (self.t_step + 1) % UPDATE_EVERY
34         if self.t_step == 0:
35
36             self.qnetwork_target.load_state_dict(self.qnetwork_local.state_dict())
37
38     def act(self, state, eps=0.):
39
40         state = torch.from_numpy(state).float().unsqueeze(0).to(device)
41         self.qnetwork_local.eval()
42         with torch.no_grad():
43             action_values = self.qnetwork_local(state)
44         self.qnetwork_local.train()
45
46         ''' Implementing Softmax Exploration Strategy'''
47         actions = action_values.cpu().data.numpy().flatten()
48
49         Qmax = np.max(actions) # For Numerical Stability
50         # To prevent overflow of values(e^(large value)) which would result in 'NaN' when computed probabilities
51         actions = actions - Qmax
52         expQ = np.exp(actions / 0.3)
53         probQ = expQ / np.sum(expQ) # Softmax Probabilities
54
55         action = np.random.choice(np.arange(self.action_size), p=probQ)
56
57         return action # Action selected based on computed softmax probabilities
58
59     def learn(self, experiences, gamma):
60         """ +E EXPERIENCE REPLAY PRESENT """
61         states, actions, rewards, next_states, dones = experiences
62
63         ''' Get max predicted Q values (for next states) from target model'''
64         Q_targets_next = self.qnetwork_target(next_states).detach().max(1)[0].unsqueeze(1)
65
66         ''' Compute Q targets for current states '''
67         Q_targets = rewards + (gamma * Q_targets_next * (1 - dones))
68
69         ''' Get expected Q values from local model '''
70         Q_expected = self.qnetwork_local(states).gather(1, actions)
71
72         ''' Compute loss '''
73         loss = F.mse_loss(Q_expected, Q_targets)
74
75         ''' Minimize the loss '''
76         self.optimizer.zero_grad()
77         loss.backward()
78
79         ''' Gradient Clipping '''
80         """ +T TRUNCATION PRESENT """
81         for param in self.qnetwork_local.parameters():
82             param.grad.data.clamp_(-1, 1)
83
84         self.optimizer.step()

```

In [9]:

```

1  ''' Defining DQN Algorithm '''
2
3  state_shape = env.observation_space.shape[0]
4  action_shape = env.action_space.n
5
6
7  def dqn(n_episodes=10000, max_t=1000, eps_start=1.0, eps_end=0.01, eps_decay=0.995):
8
9      scores_window = deque(maxlen=100)
10     ''' last 100 scores for checking if the avg is more than 195 '''
11
12     eps = eps_start
13     ''' initialize epsilon '''
14
15     for i_episode in range(1, n_episodes+1):
16         state, _ = env.reset()
17         score = 0
18         for t in range(max_t):
19             action = agent.act(state, eps)
20             next_state, reward, done, _, _ = env.step(action)
21             agent.step(state, action, reward, next_state, done)
22             state = next_state
23             score += reward
24             if done:
25                 break
26
27         scores_window.append(score)
28
29         eps = max(eps_end, eps_decay*eps)
30         ''' decrease epsilon '''
31
32         print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)), end='')
33
34         if i_episode % 100 == 0:
35             print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)))
36         if np.mean(scores_window) >= 195.0:
37             print('\nEnvironment solved in {:d} episodes! \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores
38             break
39     return True
40
41     ''' Trial run to check if algorithm runs and saves the data '''
42
43     begin_time = datetime.datetime.now()
44
45     print("Tutorial Agent with SoftMax Policy")
46     agent = TutorialAgent_SoftMax(state_size=state_shape, action_size = action_shape, seed = 0)
47     dqn()
48
49     time_taken = datetime.datetime.now() - begin_time
50
51     print(time_taken)

```

Tutorial Agent with SoftMax Policy
 Episode 64 Average Score: 204.14
 Environment solved in 64 episodes! Average Score: 204.14
 0:00:58.190546

Task 3 Solution

Modifying the DQN to return the rewards and episodes value

```
In [15]: 1  ''' Defining DQN Algorithm '''
2
3  state_shape = env.observation_space.shape[0]
4  action_shape = env.action_space.n
5
6
7  def dqn(n_episodes=10000, max_t=1000, eps_start=1.0, eps_end=0.01, eps_decay=0.995):
8
9      scores_window = deque(maxlen=100)
10     ''' last 100 scores for checking if the avg is more than 195 '''
11
12     reward_value = []
13     episode_num = []
14
15     eps = eps_start
16     ''' initialize epsilon '''
17
18     for i_episode in range(1, n_episodes+1):
19         state, _ = env.reset()
20         score = 0
21         for t in range(max_t):
22             action = agent.act(state, eps)
23             next_state, reward, done, _, _ = env.step(action)
24             agent.step(state, action, reward, next_state, done)
25             state = next_state
26             score += reward
27             if done:
28                 break
29
30         scores_window.append(score)
31         reward_value.append(np.mean(scores_window)) # Stores the Avg reward value
32         episode_num.append(i_episode) # To Store the episode value
33
34         eps = max(eps_end, eps_decay*eps)
35         ''' decrease epsilon '''
36
37         print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)), end="")
38
39         if i_episode % 100 == 0:
40             print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)))
41         if np.mean(scores_window) >= 195.0:
42             print('\nEnvironment solved in {:d} episodes! \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores
43             break
44     return True, reward_value, episode_num
45
```

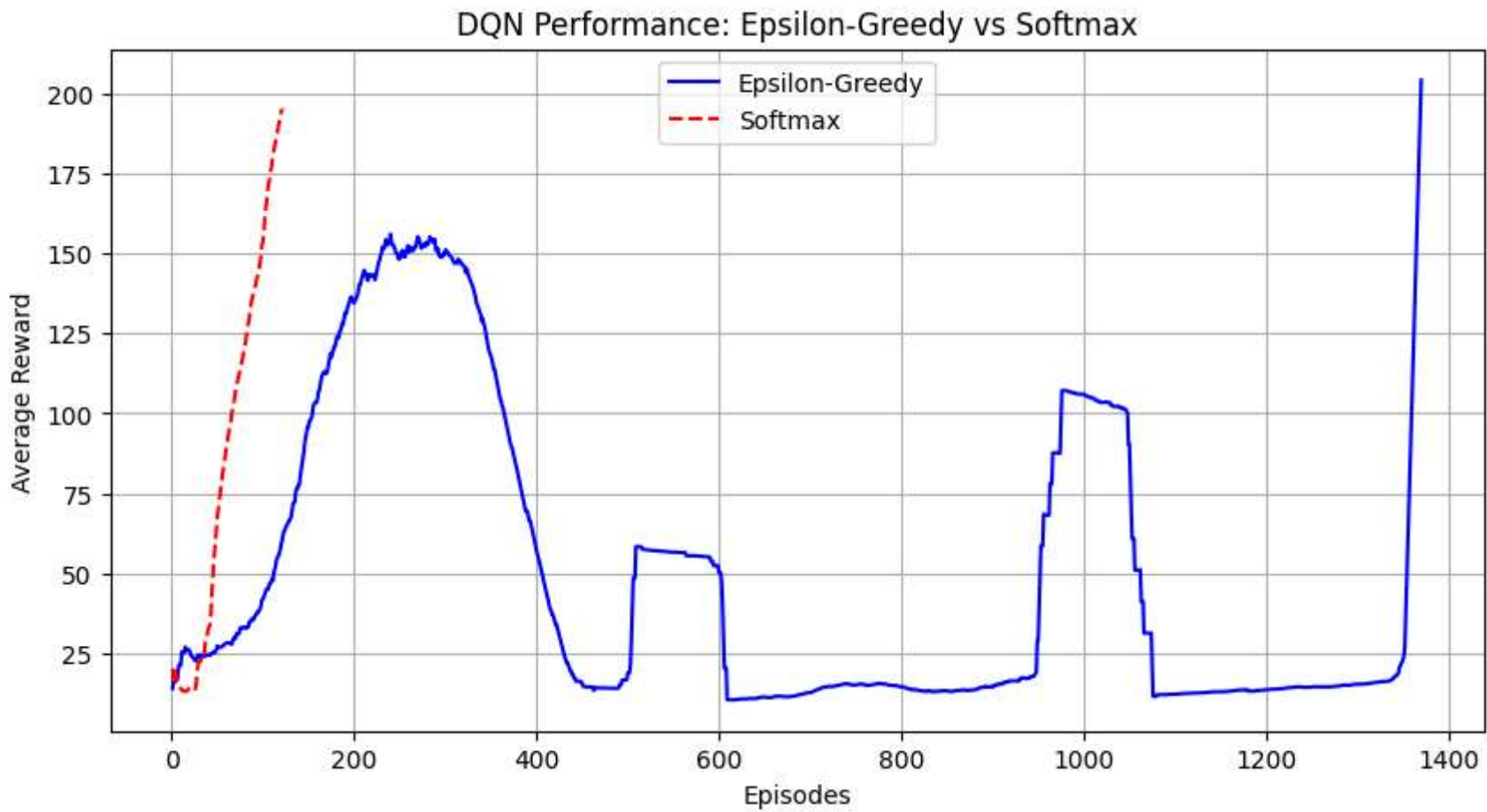
In [20]:

```
1 print("Tutorial Agent with Epsilon Greedy Policy")
2 begin_time = datetime.datetime.now()
3
4 agent = TutorialAgent(state_size=state_shape,action_size = action_shape,seed = 0)
5 _, rewards_eps, episodes_eps = dqn()
6 #print("Rewards: ", rewards_eps)
7 #print("Episodes: ", episodes_eps)
8 time_taken_eps = datetime.datetime.now() - begin_time
9
10 print("Time taken by Agent with Epsilon Greedy Policy: ", time_taken_eps)
11
12 print("Tutorial Agent with SoftMax Policy")
13 begin_time = datetime.datetime.now()
14
15 agent = TutorialAgent_SoftMax(state_size=state_shape,action_size = action_shape,seed = 0)
16 _, rewards_softmax, episodes_softmax = dqn()
17 #print("Rewards: ", rewards_softmax)
18 #print("Episodes: ", episodes_softmax)
19 time_taken_softmax = datetime.datetime.now() - begin_time
20
21 print("Time taken by Agent with Softmax Policy: ", time_taken_softmax)
```

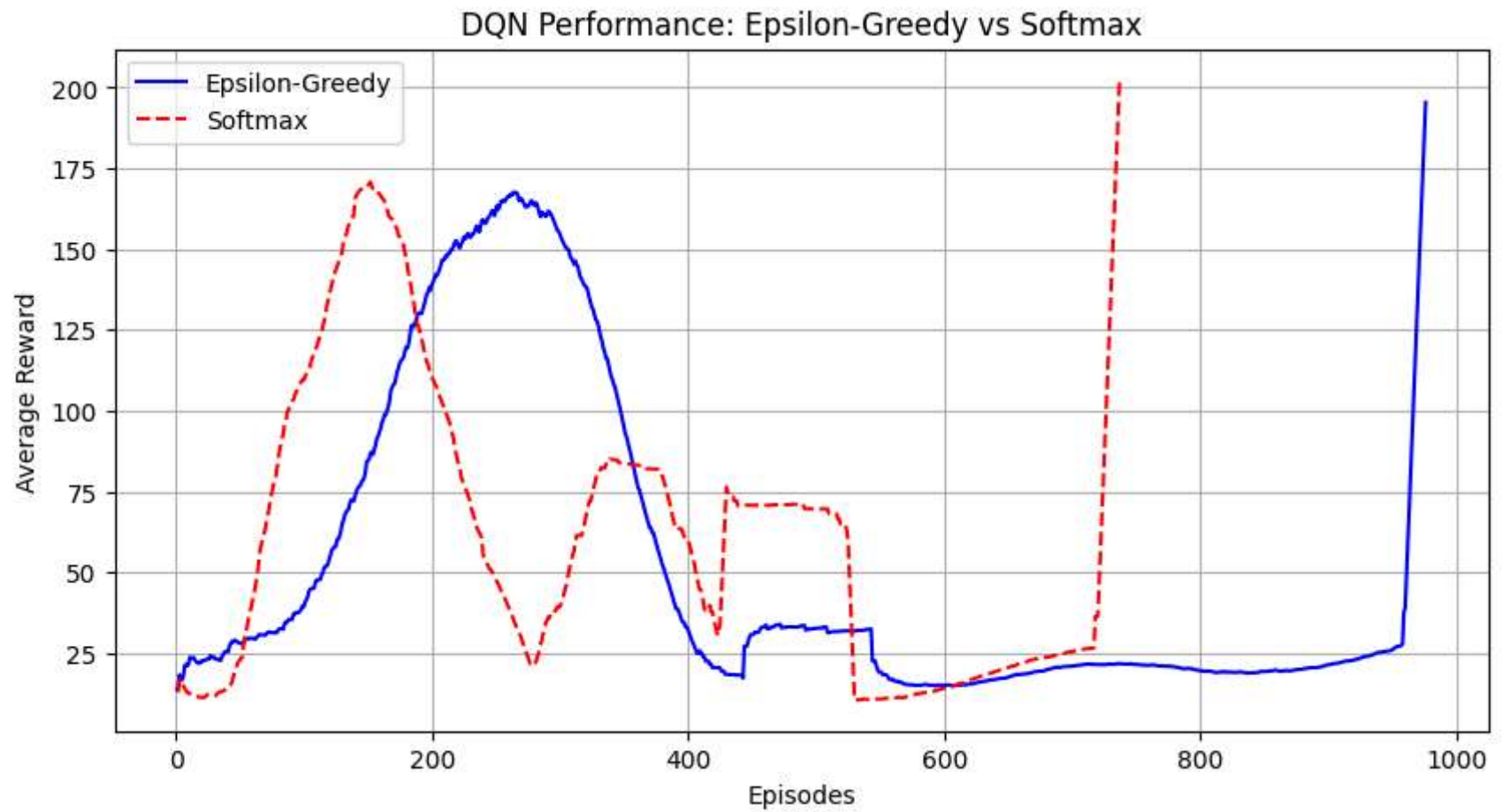
Tutorial Agent with Epsilon Greedy Policy
Episode 100 Average Score: 38.98
Episode 200 Average Score: 128.71
Episode 300 Average Score: 52.69
Episode 400 Average Score: 21.55
Episode 500 Average Score: 56.98
Episode 600 Average Score: 13.62
Episode 700 Average Score: 37.22
Episode 800 Average Score: 25.40
Episode 900 Average Score: 61.43
Episode 1000 Average Score: 87.57
Episode 1100 Average Score: 155.28
Episode 1200 Average Score: 158.14
Episode 1300 Average Score: 179.65
Episode 1400 Average Score: 176.57
Episode 1423 Average Score: 201.02
Environment solved in 1423 episodes! Average Score: 201.02
Time taken by Agent with Epsilon Greedy Policy: 0:08:19.627216
Tutorial Agent with SoftMax Policy
Episode 100 Average Score: 171.35
Episode 107 Average Score: 204.01
Environment solved in 107 episodes! Average Score: 204.01
Time taken by Agent with Softmax Policy: 0:01:20.224017

In [17]:

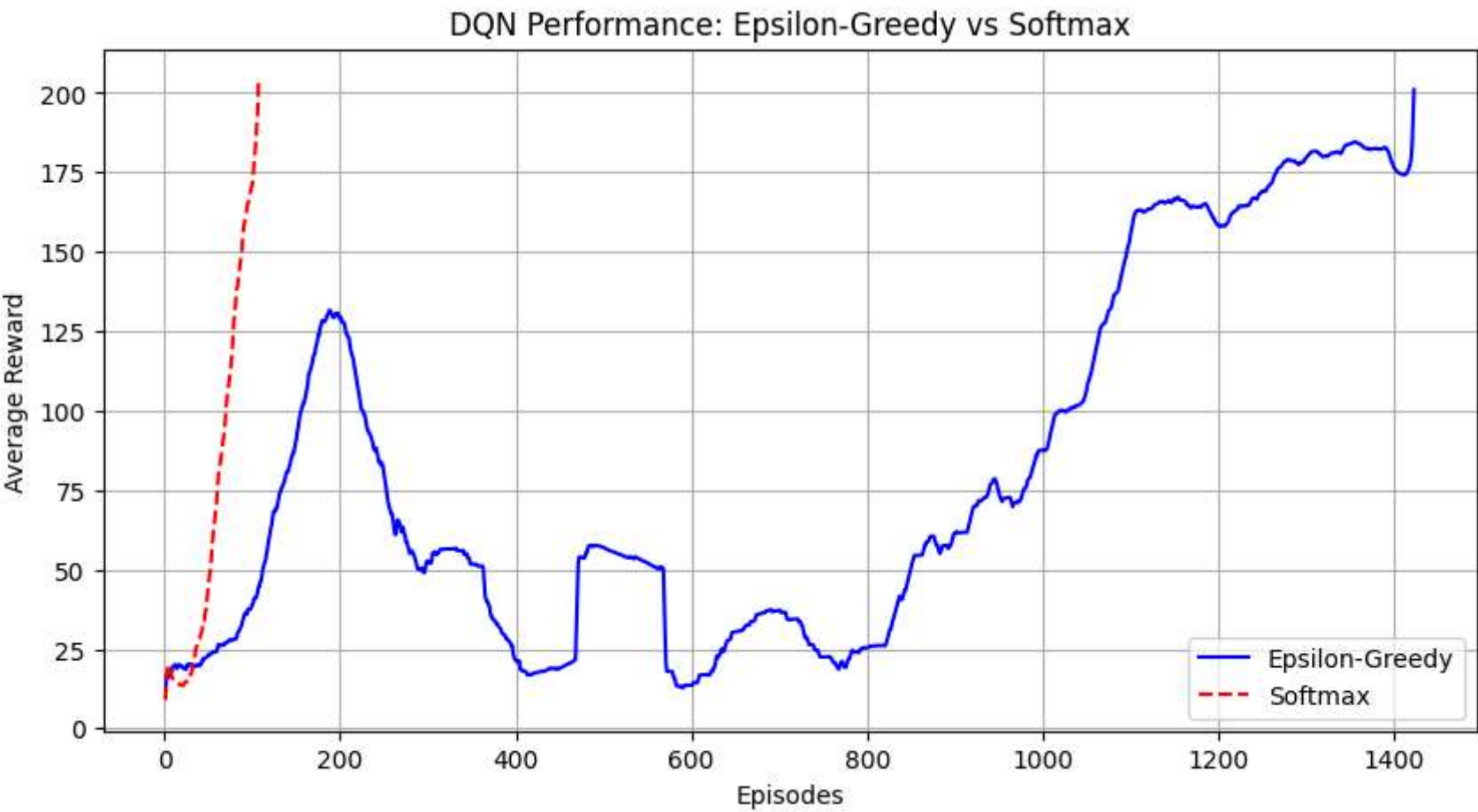
```
1 import matplotlib.pyplot as plt
2
3 # Trail 1
4 plt.figure(figsize=(10, 5))
5
6 plt.plot(epsisodes_eps, rewards_eps, label="Epsilon-Greedy", color='blue', linestyle='-')
7 plt.plot(epsisodes_softmax, rewards_softmax, label="Softmax", color='red', linestyle='--')
8
9 plt.xlabel("Episodes")
10 plt.ylabel("Average Reward")
11 plt.title("DQN Performance: Epsilon-Greedy vs Softmax")
12 plt.legend()
13 plt.grid()
14
15 plt.show()
16
```



```
In [19]: 1 import matplotlib.pyplot as plt
2
3 # Trail 2
4 plt.figure(figsize=(10, 5))
5
6 plt.plot(epsisodes_eps, rewards_eps, label="Epsilon-Greedy", color='blue', linestyle='-')
7 plt.plot(epsisodes_softmax, rewards_softmax, label="Softmax", color='red', linestyle='--')
8
9 plt.xlabel("Episodes")
10 plt.ylabel("Average Reward")
11 plt.title("DQN Performance: Epsilon-Greedy vs Softmax")
12 plt.legend()
13 plt.grid()
14
15 plt.show()
16
```




```
In [21]: 1 import matplotlib.pyplot as plt
2
3 # Trail 3
4 plt.figure(figsize=(10, 5))
5
6 plt.plot(episodes_eps, rewards_eps, label="Epsilon-Greedy", color='blue', linestyle='-')
7 plt.plot(episodes_softmax, rewards_softmax, label="Softmax", color='red', linestyle='--')
8
9 plt.xlabel("Episodes")
10 plt.ylabel("Average Reward")
11 plt.title("DQN Performance: Epsilon-Greedy vs Softmax")
12 plt.legend()
13 plt.grid()
14
15 plt.show()
16
```



DQN with a softmax strategy converges faster because it balances exploration and exploitation more effectively, leading to quicker identification of optimal actions. The controlled randomness from the temperature parameter ($\tau=0.25-0.30$) ensures steady learning without excessive exploration, improving stability.

The reward plots, compute time and the number of episode confirm the same that the DQN with softmax exploration strategy outperforms the DQN algorithm with Epsilon greedy strategy.

Epsilon Greedy Strategy

Episodes	Compute Time (min:sec)
1378	7 : 55
976	4 : 14
1423	8 : 19

Softmax Strategy

Episodes	Compute Time (min:sec)
121	1 : 45
737	3 : 58
107	1 : 20

The temperature parameter of the softmax is found by hypertuning, by running the DQN with different values of τ .