

## ▼ Tutorial 5 - DQN and Actor-Critic

Please follow this tutorial to understand the structure (code) of DQNs & get familiar with Actor Critic methods.

### References:

Please follow [Human-level control through deep reinforcement learning](#) for the original publication as well as the psuedocode. Watch Prof. Ravi's lectures on moodle or npTEL for further understanding the core concepts. Contact the TAs for further resources if needed.

### Part 1: DQN

```
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]
```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>  
Requirement already satisfied: setuptools in /usr/local/lib/python3.8/dist-packages (57.4.0)  
Collecting setuptools  
 Downloading setuptools-67.4.0-py3-none-any.whl (1.1 MB)  
 1.1/1.1 MB 37.5 MB/s eta 0:00:00  
Installing collected packages: setuptools  
Attempting uninstall: setuptools  
 Found existing installation: setuptools 57.4.0  
 Uninstalling setuptools-57.4.0:  
 Successfully uninstalled setuptools-57.4.0  
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the so  
ipython 7.9.0 requires jedi>=0.10, which is not installed.  
cvxpy 1.2.3 requires setuptools<=64.0.2, but you have setuptools 67.4.0 which is incompatible.  
Successfully installed setuptools-67.4.0  
Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>  
Requirement already satisfied: gym[classic\_control] in /usr/local/lib/python3.8/dist-packages (0.25.2)  
Requirement already satisfied: gym-notices>=0.0.4 in /usr/local/lib/python3.8/dist-packages (from gym[classic\_control]) (0.0.8)  
Requirement already satisfied: numpy>=1.18.0 in /usr/local/lib/python3.8/dist-packages (from gym[classic\_control]) (1.22.4)  
Requirement already satisfied: importlib-metadata>=4.8.0 in /usr/local/lib/python3.8/dist-packages (from gym[classic\_control]) (6.0)  
Requirement already satisfied: cloudpickle>=1.2.0 in /usr/local/lib/python3.8/dist-packages (from gym[classic\_control]) (2.2.1)  
Collecting pygame==2.1.0  
 Downloading pygame-2.1.0-cp38-cp38-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl (18.3 MB)  
 18.3/18.3 MB 74.4 MB/s eta 0:00:00  
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.8/dist-packages (from importlib-metadata>=4.8.0->gym[classic\_con  
Installing collected packages: pygame  
Successfully installed pygame-2.1.0

```
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 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
```

```

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
15 env = gym.make('CartPole-v1')
16 env.seed(0)
17
18 state_shape = env.observation_space.shape[0]
19 no_of_actions = env.action_space.n
20
21 print(state_shape)
22 print(no_of_actions)
23 print(env.action_space.sample())
24 print("----")
25
26 '''
27 # Understanding State, Action, Reward Dynamics
28
29 The agent decides an action to take depending on the state.
30
31 The Environment keeps a variable specifically for the current state.
32 - Everytime an action is passed to the environment, it calculates the new state and updates the current state variable.
33 - It returns the new current state and reward for the agent to take the next action
34
35 '''
36
37 state = env.reset()
38 ''' This returns the initial state (when environment is reset) '''
39
40 print(state)
41 print("----")
42
43 action = env.action_space.sample()
44 ''' We take a random action now '''
45
46 print(action)
47 print("----")
48
49 next_state, reward, done, info = env.step(action)
50 ''' env.step is used to calculate new state and obtain reward based on old state and action taken '''
51
52 print(next_state)
53 print(reward)
54 print(done)
55 print(info)
56 print("----")
57
4
2
0
----
[ 0.01369617 -0.02302133 -0.04590265 -0.04834723]
----
1
----
[ 0.01323574  0.17272775 -0.04686959 -0.3551522 ]
1.0
False
{}
----
/usr/local/lib/python3.8/dist-packages/gym/core.py:317: DeprecationWarning: WARN: Initializing wrapper in old step API which return
deprecation(
/usr/local/lib/python3.8/dist-packages/gym/wrappers/step_api_compatibility.py:39: DeprecationWarning: WARN: Initializing environmen
deprecation(
/usr/local/lib/python3.8/dist-packages/gym/core.py:256: DeprecationWarning: WARN: Function `env.seed(seed)` is marked as deprecated
deprecation(

```

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

```
1 '''
2 ### Q Network & Some 'hyperparameters'
3
4 QNetwork1:
5 Input Layer - 4 nodes (State Shape) \
6 Hidden Layer 1 - 64 nodes \
7 Hidden Layer 2 - 64 nodes \
8 Output Layer - 2 nodes (Action Space) \
9 Optimizer - zero_grad()
10
11 QNetwork2: Feel free to experiment more
12 '''
13
14 import torch
15 import torch.nn as nn
16 import torch.nn.functional as F
17
18
19 '''
20 Bunch of Hyper parameters (Which you might have to tune later **wink wink**)
21 '''
22 BUFFER_SIZE = int(1e5) # replay buffer size
23 BATCH_SIZE = 64 # minibatch size
24 GAMMA = 0.99 # discount factor
25 LR = 5e-4 # learning rate
26 UPDATE_EVERY = 20 # how often to update the network (When Q target is present)
27
28
29 class QNetwork1(nn.Module):
30
31     def __init__(self, state_size, action_size, seed, fc1_units=128, fc2_units=64):
32         """Initialize parameters and build model.
33         Params
34         =====
35             state_size (int): Dimension of each state
36             action_size (int): Dimension of each action
37             seed (int): Random seed
38             fc1_units (int): Number of nodes in first hidden layer
39             fc2_units (int): Number of nodes in second hidden layer
40         """
41         super(QNetwork1, self).__init__()
42         self.seed = torch.manual_seed(seed)
43         self.fc1 = nn.Linear(state_size, fc1_units)
44         self.fc2 = nn.Linear(fc1_units, fc2_units)
45         self.fc3 = nn.Linear(fc2_units, action_size)
46
47     def forward(self, state):
48         """Build a network that maps state -> action values."""
49         x = F.relu(self.fc1(state))
50         x = F.relu(self.fc2(x))
51         return self.fc3(x)
```

## ▼ Replay Buffer:

This is a 'deque' that helps us store experiences. Recall why we use such a technique.

```
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         Params
14         =====
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)

```

## ▼ Truncation:

We add a line (optionally) in the code to truncate the gradient in hopes that it would help with the stability of the learning process.

## Tutorial Agent Code:

```

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.

```

1 ''' Defining DQN Algorithm '''
2
3 state_shape = env.observation_space.shape[0]
4 action_shape = env.action_space.n
5
6 def dqn(n_episodes=10000, max_t=1000, eps_start=1.0, eps_end=0.01, eps_decay=0.995):
7
8     scores = []
9     ''' list containing scores from each episode '''
10
11     scores_window_printing = deque(maxlen=10)
12     ''' For printing in the graph '''
13
14     scores_window = deque(maxlen=100)
15     ''' last 100 scores for checking if the avg is more than 195 '''
16
17     eps = eps_start
18     ''' initialize epsilon '''
19
20     for i_episode in range(1, n_episodes+1):
21         state = env.reset()
22         score = 0
23         for t in range(max_t):
24             action = agent.act(state, eps)
25             next_state, reward, done, _ = env.step(action)
26             agent.step(state, action, reward, next_state, done)
27             state = next_state
28             score += reward
29             if done:
30                 break
31
32         scores_window.append(score)
33         scores_window_printing.append(score)
34         ''' save most recent score '''
35
36         eps = max(eps_end, eps_decay*eps)
37         ''' decrease epsilon '''
38
39         print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)), end="")
40         if i_episode % 10 == 0:
41             scores.append(np.mean(scores_window_printing))

```

```

42         if i_episode % 100 == 0:
43             print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)))
44         if np.mean(scores_window) >= 195.0:
45             print('\nEnvironment solved in {:d} episodes! \tAverage Score: {:.2f}'.format(i_episode-100, np.mean(scores_window)))
46             break
47         return [np.array(scores), i_episode-100]
48
49 ''' Trial run to check if algorithm runs and saves the data '''
50
51 begin_time = datetime.datetime.now()
52 agent = TutorialAgent(state_size=state_shape, action_size = action_shape, seed = 0)
53
54
55 dqn()
56
57
58 time_taken = datetime.datetime.now() - begin_time
59
60 print(time taken)
    Episode 100      Average Score: 38.24
    Episode 200      Average Score: 144.32
    Episode 231      Average Score: 195.80
    Environment solved in 131 episodes!      Average Score: 195.80
    0:01:57.681861

```

## ▼ 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 ( $\epsilon$ -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  $\epsilon$ -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 1a

Epsilon Greedy exploration strategy was used.

## ▼ Task 1b

Out of the two exploration strategies discussed in class ( $\epsilon$ -greedy & Softmax),  $\epsilon$ -greedy was implemented already. Implemented the strategy that's not used here i.e Softmax Policy.

```

1 from scipy.special import softmax
2 class TutorialAgent_Softmax():
3
4     def __init__(self, state_size, action_size, seed):
5
6         ''' Agent Environment Interaction '''
7         self.state_size = state_size
8         self.action_size = action_size
9         self.seed = random.seed(seed)
10
11         ''' Q-Network '''

```

```

12     self.qnetwork_local = QNetwork1(state_size, action_size, seed).to(device)
13     self.qnetwork_target = QNetwork1(state_size, action_size, seed).to(device)
14     self.optimizer = optim.Adam(self.qnetwork_local.parameters(), lr=LR)
15
16     ''' Replay memory '''
17     self.memory = ReplayBuffer(action_size, BUFFER_SIZE, BATCH_SIZE, seed)
18
19     ''' Initialize time step (for updating every UPDATE_EVERY steps)          -Needed for Q Targets '''
20     self.t_step = 0
21
22 def step(self, state, action, reward, next_state, done):
23
24     ''' Save experience in replay memory '''
25     self.memory.add(state, action, reward, next_state, done)
26
27     ''' If enough samples are available in memory, get random subset and learn '''
28     if len(self.memory) >= BATCH_SIZE:
29         experiences = self.memory.sample()
30         self.learn(experiences, GAMMA)
31
32     """ +Q TARGETS PRESENT """
33     ''' Updating the Network every 'UPDATE_EVERY' steps taken '''
34     self.t_step = (self.t_step + 1) % UPDATE_EVERY
35     if self.t_step == 0:
36
37         self.qnetwork_target.load_state_dict(self.qnetwork_local.state_dict())
38
39 def act(self, state, temp = 1):
40
41     state = torch.from_numpy(state).float().unsqueeze(0).to(device)
42     self.qnetwork_local.eval()
43     with torch.no_grad():
44         action_values = self.qnetwork_local(state)
45     self.qnetwork_local.train()
46
47     ''' softmax action selection '''
48     prob = np.nan_to_num(torch.softmax(action_values.cpu().data.numpy()[0]/temp))
49     prob /= prob.sum()
50     return np.random.choice(self.action_size, p = prob)
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()
78
79
80 1 ''' Defining DQN Algorithm '''
81 2
82 3 state_shape = env.observation_space.shape[0]
83 4 action_shape = env.action_space.n
84 5
85 6 def dqn_softmax(n_episodes=10000, max_t=1000, temp_start= 10000, temp_end=1, temp_decay=0.995):
86 7
87 8     scores_softmax = []
88 9     ''' list containing scores from each episode '''
89
90 10
91 11     scores_window_printing = deque(maxlen=10)
92 12     ''' For printing in the graph '''
93 13
94 14     scores_window= deque(maxlen=100)

```

```

15     ''' last 100 scores for checking if the avg is more than 195 '''
16
17     temp = temp_start
18     ''' initialize epsilon '''
19
20     for i_episode in range(1, n_episodes+1):
21         state = env.reset()
22         score = 0
23         for t in range(max_t):
24             action = agent.act(state, temp)
25             #print(action)
26             next_state, reward, done, _ = env.step(action)
27             agent.step(state, action, reward, next_state, done)
28             state = next_state
29             score += reward
30             if done:
31                 break
32
33         scores_window.append(score)
34         scores_window_printing.append(score)
35         ''' save most recent score '''
36
37         temp = temp*temp_decay
38         ''' decrease temperature '''
39
40         print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)), end="")
41
42         scores_softmax.append(np.mean(scores_window_printing))
43         if i_episode % 100 == 0:
44             print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_window)))
45         if np.mean(scores_window) >= 195.0:
46             print('\nEnvironment solved in {:d} episodes! \tAverage Score: {:.2f}'.format(i_episode-100, np.mean(scores_window)))
47             break
48     return [np.array(scores_softmax), i_episode-100]
49
50 ''' Trial run to check if algorithm runs and saves the data '''
51
52 begin_time = datetime.datetime.now()
53 agent = TutorialAgent_Softmax(state_size=state_shape, action_size = action_shape, seed = 0)
54
55
56 dqn_softmax()
57
58
59 time_taken = datetime.datetime.now() - begin_time
60
61 print(time_taken)

```

```

Episode 100      Average Score: 24.19
Episode 200      Average Score: 20.25
Episode 300      Average Score: 23.88
Episode 400      Average Score: 24.72
Episode 500      Average Score: 28.62
Episode 600      Average Score: 35.70
Episode 700      Average Score: 100.10
Episode 787      Average Score: 195.21
Environment solved in 687 episodes!      Average Score: 195.21
0:02:53.279581

```

## ▼ 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  $\epsilon$ -greedy) with DQN (using softmax). Think along the lines of 'no. of episodes', 'reward plots', 'compute time', etc. and add a few comments.

Computation Time =>  $\epsilon$ -greedy = 01:54 , softmax = 02:43

No. of episodes till convergence =>  $\epsilon$ -greedy = 131 , softmax = 658

Softmax is slower than  $\epsilon$ -greedy in terms of "compute time" and takes more "number of episodes" to converge.

```

1  #running with epsilon-greedy strategy
2  begin_time = datetime.datetime.now()
3  agent = TutorialAgent(state_size=state_shape, action_size = action_shape, seed = 0)

```



```

4
5 [y_axis, x_axis] = dqn()
6
7 time_taken = datetime.datetime.now() - begin_time
8
9 print("Computation Time =" , time_taken)

Episode 100      Average Score: 42.36
Episode 200      Average Score: 125.91
Episode 236      Average Score: 196.09
Environment solved in 136 episodes!      Average Score: 196.09
Computation Time = 0:01:38.273009

1 scores_eps_greedy = []
2 N_ep = len(y_axis)
3 for i in range(N_ep):
4     idx = min(99,i)
5     scores_eps_greedy.append(np.mean(y_axis[i-idx:i+1]))

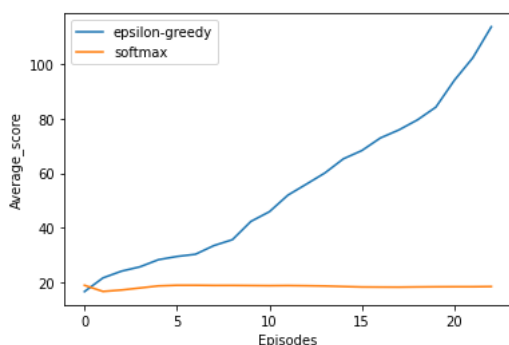
1 begin_time = datetime.datetime.now()
2 agent = TutorialAgent_Softmax(state_size=state_shape,action_size = action_shape,seed = 0)
3
4 [y_axis_softmax ,x_axis_softmax] = dqn_softmax()
5
6 time_taken_softmax = datetime.datetime.now() - begin_time
7
8 print(time_taken_softmax)

Episode 100      Average Score: 22.55
Episode 200      Average Score: 22.96
Episode 300      Average Score: 20.76
Episode 400      Average Score: 21.48
Episode 500      Average Score: 23.51
Episode 600      Average Score: 28.74
Episode 700      Average Score: 40.66
Episode 800      Average Score: 135.12
Episode 900      Average Score: 184.15
Episode 1000     Average Score: 171.38
Episode 1100     Average Score: 155.53
Episode 1200     Average Score: 104.70
Episode 1300     Average Score: 56.85
Episode 1400     Average Score: 34.08
Episode 1500     Average Score: 22.37
Episode 1600     Average Score: 42.32
Episode 1700     Average Score: 39.83
Episode 1800     Average Score: 158.11
Episode 1900     Average Score: 34.87
Episode 1987     Average Score: 195.24
Environment solved in 1887 episodes!      Average Score: 195.24
0:10:15.072345

1 scores_softmax = []
2 N_ep = len(y_axis)
3 for i in range(N_ep):
4     idx = min(99,i)
5     scores_softmax.append(np.mean(y_axis_softmax[i-idx:i+1]))

1 ### Plot of total reward vs episode
2 ## Write Code Below
3
4 plt.xlabel('Episodes')
5 plt.ylabel('Average_score')
6 plt.plot(np.arange(len(scores_eps_greedy)) ,scores_eps_greedy)
7 plt.plot(np.arange(len(scores_softmax)) ,scores_softmax)
8 plt.legend(['epsilon-greedy', 'softmax'])
9 plt.show()

```



## Part 2: One-Step Actor-Critic Algorithm

**Actor-Critic methods** learn both a policy  $\pi(a|s; \theta)$  and a state-value function  $v(s; w)$  simultaneously. The policy is referred to as the actor that suggests actions given a state. The estimated value function is referred to as the critic. It evaluates actions taken by the actor based on the given policy. In this exercise, both functions are approximated by feedforward neural networks.

- The policy network is parametrized by  $\theta$  - it takes a state  $s$  as input and outputs the probabilities  $\pi(a|s; \theta) \forall a$
- The value network is parametrized by  $w$  - it takes a state  $s$  as input and outputs a scalar value associated with the state, i.e.,  $v(s; w)$
- The single step TD error can be defined as follows:

$$\delta_t = R_{t+1} + \gamma v(s_{t+1}; w) - v(s_t; w)$$

- The loss function to be minimized at every step ( $L_{tot}^{(t)}$ ) is a summation of two terms, as follows:

$$L_{tot}^{(t)} = L_{actor}^{(t)} + L_{critic}^{(t)}$$

where,

$$L_{actor}^{(t)} = -\log \pi(a_t|s_t; \theta) \delta_t$$

$$L_{critic}^{(t)} = \delta_t^2$$

- **NOTE: Here, weights of the first two hidden layers are shared by the policy and the value network**
  - First two hidden layer sizes: [1024, 512]
  - Output size of policy network: 2 (Softmax activation)
  - Output size of value network: 1 (Linear activation)

## Initializing Actor-Critic Network

```
1 class ActorCriticModel(tf.keras.Model):
2     """
3     Defining policy and value networkss
4     """
5     def __init__(self, action_size, n_hidden1=1024, n_hidden2=512):
6         super(ActorCriticModel, self).__init__()
7
8         #Hidden Layer 1
9         self.fc1 = tf.keras.layers.Dense(n_hidden1, activation='relu')
10        #Hidden Layer 2
11        self.fc2 = tf.keras.layers.Dense(n_hidden2, activation='relu')
12
13        #Output Layer for policy
14        self.pi_out = tf.keras.layers.Dense(action_size, activation='softmax')
15        #Output Layer for state-value
16        self.v_out = tf.keras.layers.Dense(1)
17
18    def call(self, state):
19        """
20        Computes policy distribution and state-value for a given state
21        """
22        layer1 = self.fc1(state)
23        layer2 = self.fc2(layer1)
24
25        pi = self.pi_out(layer2)
26        v = self.v_out(layer2)
27
28        return pi, v
```

## Agent Class

**Task 2a:** Write code to compute  $\delta_t$  inside the Agent.learn() function

```
1 class Agent:
2     """
3     Agent class
4     """
5     def __init__(self, action_size, lr=0.001, gamma=0.99, seed = 85):
6         self.gamma = gamma
7         self.ac_model = ActorCriticModel(action_size=action_size)
8         self.ac_model.compile(tf.keras.optimizers.Adam(learning_rate=lr))
9         np.random.seed(seed)
10
11    def sample_action(self, state):
12        """
13        Given a state, compute the policy distribution over all actions and sample one action
14        """
```

```

15         pi,_ = self.ac_model(state)
16
17         action_probabilities = tfp.distributions.Categorical(probs=pi)
18         sample = action_probabilities.sample()
19
20         return int(sample.numpy()[0])
21
22     def actor_loss(self, action, pi, delta):
23         """
24         Compute Actor Loss
25         """
26         return -tf.math.log(pi[0,action]) * delta
27
28     def critic_loss(self,delta):
29         """
30         Critic loss aims to minimize TD error
31         """
32         return delta**2
33
34     @tf.function
35     def learn(self, state, action, reward, next_state, done):
36         """
37         For a given transition (s,a,s',r) update the paramters by computing the
38         gradient of the total loss
39         """
40         with tf.GradientTape(persistent=True) as tape:
41             pi, V_s = self.ac_model(state)
42             _, V_s_next = self.ac_model(next_state)
43
44             V_s = tf.squeeze(V_s)
45             V_s_next = tf.squeeze(V_s_next)
46
47
48             ##### TO DO: Write the equation for delta (TD error)
49             ## Write code below
50             delta = delta = reward + self.gamma * V_s_next * (1 - done) - V_s ## Complete this
51             loss_a = self.actor_loss(action, pi, delta)
52             loss_c =self.critic_loss(delta)
53             loss_total = loss_a + loss_c
54
55             gradient = tape.gradient(loss_total, self.ac_model.trainable_variables)
56             self.ac_model.optimizer.apply_gradients(zip(gradient, self.ac_model.trainable_variables))

```

## ▼ Train the Network

```

1 env = gym.make('CartPole-v1')
2
3 #Initializing Agent
4 agent = Agent(lr=1e-4, action_size=env.action_space.n)
5 #Number of episodes
6 episodes = 1800
7 tf.compat.v1.reset_default_graph()
8
9 reward_list = []
10 average_reward_list = []
11 begin_time = datetime.datetime.now()
12
13 for ep in range(1, episodes + 1):
14     state = env.reset().reshape(1,-1)
15     done = False
16     ep_rew = 0
17     while not done:
18         action = agent.sample_action(state) ##Sample Action
19         next_state, reward, done, info = env.step(action) ##Take action
20         next_state = next_state.reshape(1,-1)
21         ep_rew += reward ##Updating episode reward
22         agent.learn(state, action, reward, next_state, done) ##Update Parameters
23         state = next_state ##Updating State
24     reward_list.append(ep_rew)
25
26     if ep % 10 == 0:
27         avg_rew = np.mean(reward_list[-10:])
28         print('Episode ', ep, 'Reward %f' % ep_rew, 'Average Reward %f' % avg_rew)
29
30     if ep % 100:
31         avg_100 = np.mean(reward_list[-100:])
32         if avg_100 > 195.0:
33             print('Stopped at Episode ',ep-100)
34             break
35

```

```

36 time_taken = datetime.datetime.now() - begin_time
37 print(time_taken)

Episode 10 Reward 19.000000 Average Reward 17.200000
Episode 20 Reward 13.000000 Average Reward 30.800000
Episode 30 Reward 14.000000 Average Reward 24.700000
Episode 40 Reward 11.000000 Average Reward 21.700000
Episode 50 Reward 12.000000 Average Reward 17.400000
Episode 60 Reward 9.000000 Average Reward 14.900000
Episode 70 Reward 19.000000 Average Reward 13.200000
Episode 80 Reward 11.000000 Average Reward 13.400000
Episode 90 Reward 16.000000 Average Reward 27.600000
Episode 100 Reward 18.000000 Average Reward 30.700000
Episode 110 Reward 31.000000 Average Reward 43.500000
Episode 120 Reward 127.000000 Average Reward 54.100000
Episode 130 Reward 46.000000 Average Reward 70.700000
Episode 140 Reward 51.000000 Average Reward 61.500000
Episode 150 Reward 46.000000 Average Reward 70.900000
Episode 160 Reward 103.000000 Average Reward 110.300000
Episode 170 Reward 111.000000 Average Reward 102.700000
Episode 180 Reward 121.000000 Average Reward 116.500000
Episode 190 Reward 136.000000 Average Reward 117.400000
Episode 200 Reward 107.000000 Average Reward 111.600000
Episode 210 Reward 102.000000 Average Reward 99.100000
Episode 220 Reward 121.000000 Average Reward 114.200000
Episode 230 Reward 126.000000 Average Reward 90.700000
Episode 240 Reward 103.000000 Average Reward 107.800000
Episode 250 Reward 111.000000 Average Reward 112.500000
Episode 260 Reward 128.000000 Average Reward 112.100000
Episode 270 Reward 109.000000 Average Reward 116.300000
Episode 280 Reward 58.000000 Average Reward 106.400000
Episode 290 Reward 108.000000 Average Reward 107.000000
Episode 300 Reward 123.000000 Average Reward 106.200000
Episode 310 Reward 106.000000 Average Reward 119.600000
Episode 320 Reward 107.000000 Average Reward 101.400000
Episode 330 Reward 114.000000 Average Reward 103.500000
Episode 340 Reward 134.000000 Average Reward 120.700000
Episode 350 Reward 75.000000 Average Reward 59.000000
Episode 360 Reward 119.000000 Average Reward 86.200000
Episode 370 Reward 128.000000 Average Reward 113.800000
Episode 380 Reward 120.000000 Average Reward 126.900000
Episode 390 Reward 143.000000 Average Reward 116.800000
Episode 400 Reward 159.000000 Average Reward 134.900000
Episode 410 Reward 102.000000 Average Reward 151.400000
Episode 420 Reward 248.000000 Average Reward 152.400000
Episode 430 Reward 204.000000 Average Reward 214.500000
Episode 440 Reward 138.000000 Average Reward 137.800000
Episode 450 Reward 113.000000 Average Reward 130.300000
Episode 460 Reward 123.000000 Average Reward 168.700000
Episode 470 Reward 101.000000 Average Reward 109.600000
Episode 480 Reward 122.000000 Average Reward 119.000000
Episode 490 Reward 125.000000 Average Reward 116.400000
Episode 500 Reward 123.000000 Average Reward 136.400000
Episode 510 Reward 432.000000 Average Reward 141.000000
Episode 520 Reward 293.000000 Average Reward 160.800000
Episode 530 Reward 112.000000 Average Reward 119.100000
Episode 540 Reward 160.000000 Average Reward 125.100000
Episode 550 Reward 161.000000 Average Reward 148.400000
Episode 560 Reward 110.000000 Average Reward 125.900000
Episode 570 Reward 122.000000 Average Reward 125.600000
Episode 580 Reward 93.000000 Average Reward 99.500000

```

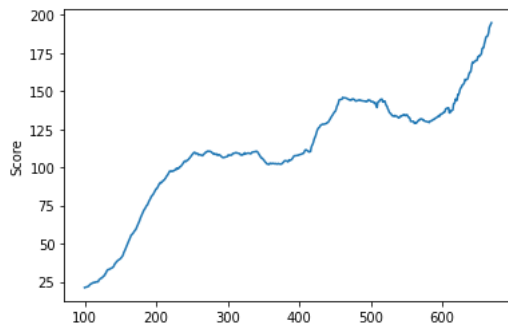
## ▼ Task 2b: Plot total reward curve

In the cell below, write code to plot the total reward averaged over 100 episodes (moving average)

```

1 ### Plot of total reward vs episode
2 ## Write Code Below
3
4 def plot_scores(scores):
5     # calculate moving average over 100 episodes
6     ma_scores = [np.mean(scores[i:i+100]) for i in range(len(scores)-100)]
7
8     # plot scores and moving average
9     fig = plt.figure()
10    ax = fig.add_subplot(111)
11    #plt.plot(np.arange(len(scores)), scores)
12    plt.plot(np.arange(100, len(scores)), ma_scores)
13    plt.ylabel('Score')
14    plt.xlabel('Episode #')
15    plt.show()
16
17 plot_scores(reward_list)
18

```



▼ Code for rendering ([source](#))

```

1 # Render an episode and save as a GIF file
2
3 display = Display(visible=0, size=(400, 300))
4 display.start()
5
6
7 def render_episode(env: gym.Env, model: tf.keras.Model, max_steps: int):
8     screen = env.render(mode='rgb_array')
9     im = Image.fromarray(screen)
10
11     images = [im]
12
13     state = tf.constant(env.reset(), dtype=tf.float32)
14     for i in range(1, max_steps + 1):
15         state = tf.expand_dims(state, 0)
16         action_probs, _ = model(state)
17         action = np.argmax(np.squeeze(action_probs))
18         state, _, done, _ = env.step(action)
19         state = tf.constant(state, dtype=tf.float32)
20
21     # Render screen every 10 steps
22     if i % 10 == 0:
23         screen = env.render(mode='rgb_array')
24         images.append(Image.fromarray(screen))
25
26     if done:
27         break
28
29     return images
30
31
32 # Save GIF image
33 images = render_episode(env, agent.ac_model, 200)
34 image_file = 'cartpole-v1.gif'
35 # loop=0: loop forever, duration=1: play each frame for 1ms
36 images[0].save(
37     image_file, save_all=True, append_images=images[1:], loop=0, duration=1)

/usr/local/lib/python3.8/dist-packages/gym/core.py:43: DeprecationWarning: WARN: The argument mode in render method is deprecated;
See here for more information: https://www.gymnasium.dev/content/api/deprecation/
deprecation(

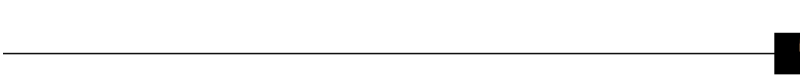
```

---

```

1 import tensorflow_docs.vis.embed as embed
2 embed.embed_file(image_file)

```



✓ 0s completed at 12:51 AM

