

Continuous_Control

August 2, 2022

1 Continuous Control

In this notebook, you will learn how to use the Unity ML-Agents environment for the second project of the [Deep Reinforcement Learning Nanodegree](#) program.

1.0.1 1. Start the Environment

We begin by importing the necessary packages. If the code cell below returns an error, please revisit the project instructions to double-check that you have installed [Unity ML-Agents](#) and [NumPy](#).

```
[1]: from unityagents import UnityEnvironment
import numpy as np
```

Next, we will start the environment! *Before running the code cell below*, change the `file_name` parameter to match the location of the Unity environment that you downloaded.

- **Mac**: "path/to/Reacher.app"
- **Windows (x86)**: "path/to/Reacher_Windows_x86/Reacher.exe"
- **Windows (x86_64)**: "path/to/Reacher_Windows_x86_64/Reacher.exe"
- **Linux (x86)**: "path/to/Reacher_Linux/Reacher.x86"
- **Linux (x86_64)**: "path/to/Reacher_Linux/Reacher.x86_64"
- **Linux (x86, headless)**: "path/to/Reacher_Linux_NoVis/Reacher.x86"
- **Linux (x86_64, headless)**: "path/to/Reacher_Linux_NoVis/Reacher.x86_64"

For instance, if you are using a Mac, then you downloaded `Reacher.app`. If this file is in the same folder as the notebook, then the line below should appear as follows:

```
env = UnityEnvironment(file_name="Reacher.app")
```

```
[2]: env = UnityEnvironment(file_name='./Reacher_Linux_NoVis/Reacher.x86_64')
```

```
INFO:unityagents:
'Academy' started successfully!
Unity Academy name: Academy
  Number of Brains: 1
  Number of External Brains : 1
  Lesson number : 0
  Reset Parameters :
    goal_speed -> 1.0
```

```

        goal_size -> 5.0
Unity brain name: ReacherBrain
    Number of Visual Observations (per agent): 0
    Vector Observation space type: continuous
    Vector Observation space size (per agent): 33
    Number of stacked Vector Observation: 1
    Vector Action space type: continuous
    Vector Action space size (per agent): 4
    Vector Action descriptions: , , ,

```

Environments contain *brains* which are responsible for deciding the actions of their associated agents. Here we check for the first brain available, and set it as the default brain we will be controlling from Python.

```

[3]: # get the default brain
brain_name = env.brain_names[0]
brain = env.brains[brain_name]

```

1.0.2 2. Examine the State and Action Spaces

In this environment, a double-jointed arm can move to target locations. A reward of +0.1 is provided for each step that the agent's hand is in the goal location. Thus, the goal of your agent is to maintain its position at the target location for as many time steps as possible.

The observation space consists of 33 variables corresponding to position, rotation, velocity, and angular velocities of the arm. Each action is a vector with four numbers, corresponding to torque applicable to two joints. Every entry in the action vector must be a number between -1 and 1.

Run the code cell below to print some information about the environment.

```

[4]: # reset the environment
env_info = env.reset(train_mode=True)[brain_name]

# number of agents
num_agents = len(env_info.agents)
print('Number of agents:', num_agents)

# size of each action
action_size = brain.vector_action_space_size
print('Size of each action:', action_size)

# examine the state space
states = env_info.vector_observations
state_size = states.shape[1]
print('There are {} agents. Each observes a state with length: {}'.
      ↪format(states.shape[0], state_size))
print('The state for the first agent looks like:', states[0])

```

```

Number of agents: 20
Size of each action: 4

```

There are 20 agents. Each observes a state with length: 33
The state for the first agent looks like: [0.00000000e+00 -4.00000000e+00
0.00000000e+00 1.00000000e+00
-0.00000000e+00 -0.00000000e+00 -4.37113883e-08 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 -1.00000000e+01 0.00000000e+00
1.00000000e+00 -0.00000000e+00 -0.00000000e+00 -4.37113883e-08
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 5.75471878e+00 -1.00000000e+00
5.55726624e+00 0.00000000e+00 1.00000000e+00 0.00000000e+00
-1.68164849e-01]

1.0.3 3. Take Random Actions in the Environment

In the next code cell, you will learn how to use the Python API to control the agent and receive feedback from the environment.

Once this cell is executed, you will watch the agent's performance, if it selects an action at random with each time step. A window should pop up that allows you to observe the agent, as it moves through the environment.

Of course, as part of the project, you'll have to change the code so that the agent is able to use its experience to gradually choose better actions when interacting with the environment!

```
[5]: env_info = env.reset(train_mode=False)[brain_name]      # reset the environment
      ↪
      states = env_info.vector_observations                  # get the current state
      ↪(for each agent)
      scores = np.zeros(num_agents)                        # initialize the score
      ↪(for each agent)
      while True:
          actions = np.random.randn(num_agents, action_size) # select an action (for
          ↪each agent)
          actions = np.clip(actions, -1, 1)                 # all actions between -1
          ↪and 1
          env_info = env.step(actions)[brain_name]          # send all actions to
          ↪the environment
          next_states = env_info.vector_observations         # get next state (for
          ↪each agent)
          rewards = env_info.rewards                        # get reward (for each
          ↪agent)
          dones = env_info.local_done                       # see if episode finished
          scores += env_info.rewards                        # update the score (for
          ↪each agent)
          states = next_states                              # roll over states to
          ↪next time step
          if np.any(dones):                                 # exit loop if episode
          ↪finished
```

```

        break
print('Total score (averaged over agents) this episode: {}'.format(np.
    ↪mean(scores)))

```

Total score (averaged over agents) this episode: 0.05899999868124724

1.0.4 4. It's Your Turn!

Now it's your turn to train your own agent to solve the environment! When training the environment, set `train_mode=True`, so that the line for resetting the environment looks like the following:

```
env_info = env.reset(train_mode=True)[brain_name]
```

1.0.5 Define Actor & Critic Models

Neural network architectures - Actor - Input layer (size 33) - FC layer (size 400) - Relu - FC layer (size 300) - Relu - Output layer (size 4) - Tanh - Critic - Input layer (size 33) - FC layer (size 404) - Relu - FC layer (size 300) - Relu - Output layer (size 1)

```

[6]: import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim

from torchsummary import summary

```

```

[7]: def hidden_init(layer):
    fan_in = layer.weight.data.size()[0]
    lim = 1. / np.sqrt(fan_in)
    return (-lim, lim)

```

```

[8]: class Actor(nn.Module):
    """Actor (Policy) Model."""

    def __init__(self, state_size, action_size, seed, fc1_units=400,
    ↪fc2_units=300):
        """Initialize parameters and build model.
        Params
        =====
            state_size (int): Dimension of each state
            action_size (int): Dimension of each action
            seed (int): Random seed
            fc1_units (int): Number of nodes in first hidden layer
            fc2_units (int): Number of nodes in second hidden layer
        """
        super(Actor, self).__init__()
        self.seed = torch.manual_seed(seed)
        self.fc1 = nn.Linear(state_size, fc1_units)

```

```

self.fc2 = nn.Linear(fc1_units, fc2_units)
self.fc3 = nn.Linear(fc2_units, action_size)
self.reset_parameters()

def reset_parameters(self):
    self.fc1.weight.data.uniform_(*hidden_init(self.fc1))
    self.fc2.weight.data.uniform_(*hidden_init(self.fc2))
    self.fc3.weight.data.uniform_(-3e-3, 3e-3)

def forward(self, state):
    """Build an actor (policy) network that maps states -> actions."""
    x = F.relu(self.fc1(state))
    x = F.relu(self.fc2(x))
    return F.tanh(self.fc3(x))

```

```

[9]: class Critic(nn.Module):
    """Critic (Value) Model."""

    def __init__(self, state_size, action_size, seed, fcs1_units=400,
    ↪fc2_units=300):
        """Initialize parameters and build model.
        Params
        =====
        state_size (int): Dimension of each state
        action_size (int): Dimension of each action
        seed (int): Random seed
        fcs1_units (int): Number of nodes in the first hidden layer
        fc2_units (int): Number of nodes in the second hidden layer
        """
        super(Critic, self).__init__()
        self.seed = torch.manual_seed(seed)
        self.fcs1 = nn.Linear(state_size, fcs1_units)
        self.fc2 = nn.Linear(fcs1_units+action_size, fc2_units)
        self.fc3 = nn.Linear(fc2_units, 1)
        self.reset_parameters()

    def reset_parameters(self):
        self.fcs1.weight.data.uniform_(*hidden_init(self.fcs1))
        self.fc2.weight.data.uniform_(*hidden_init(self.fc2))
        self.fc3.weight.data.uniform_(-3e-3, 3e-3)

    def forward(self, state, action):
        """Build a critic (value) network that maps (state, action) pairs ->
    ↪Q-values."""
        xs = F.relu(self.fcs1(state))
        x = torch.cat((xs, action), dim=1)
        x = F.relu(self.fc2(x))

```

```
return self.fc3(x)
```

1.0.6 Implement Agent

The Deep Deterministic Policy Gradients (DDPG) algorithm below successfully solves the task of this project. It learns from the provided environment without any prior knowledge of it or data labels and maximizes reward by interacting with the environment.

Algorithm 1 DDPG algorithm

Randomly initialize critic network $Q(s, a|\theta^Q)$ and actor $\mu(s|\theta^\mu)$ with weights θ^Q and θ^μ .
Initialize target network Q' and μ' with weights $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$
Initialize replay buffer R
for episode = 1, M **do**
 Initialize a random process \mathcal{N} for action exploration
 Receive initial observation state s_1
 for t = 1, T **do**
 Select action $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$ according to the current policy and exploration noise
 Execute action a_t and observe reward r_t and observe new state s_{t+1}
 Store transition (s_t, a_t, r_t, s_{t+1}) in R
 Sample a random minibatch of N transitions (s_i, a_i, r_i, s_{i+1}) from R
 Set $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$
 Update critic by minimizing the loss: $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$
 Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

Update the target networks:

$$\begin{aligned} \theta^{Q'} &\leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'} \\ \theta^{\mu'} &\leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'} \end{aligned}$$

end for
end for

(descrip-

tion courtesy of <https://arxiv.org/pdf/1509.02971.pdf>, see page 5)

Based on this implementation, adaptations were needed in the model and Noise process to adapt from one to multiple agents. Moreover, the learning rate of the critic was reduced to 1e-4 which lead to faster and more stable learning.

```
[10]: import random
import copy
from collections import namedtuple, deque
```

The DDPG agent contains the following components and configs: - A replay buffer to store memories with the size of 1e5 - Minibatch sizes of 128 - A discount factor of 0.99 for value function approximation - A soft update to blend the regular into the target network of 1e-3 - Learning rates of the actor and critic each set to 1e-4 - Noise according to the Ornstein-Uhlenbeck process with theta=0.15, sigma=0.2

```
[11]: BUFFER_SIZE = int(1e5)  # replay buffer size
      BATCH_SIZE = 128      # minibatch size
      GAMMA = 0.99          # discount factor
      TAU = 1e-3            # for soft update of target parameters
      LR_ACTOR = 1e-4       # learning rate of the actor
      LR_CRITIC = 1e-4      # learning rate of the critic
      WEIGHT_DECAY = 0      # L2 weight decay
```

```
[12]: class ReplayBuffer:
      """Fixed-size buffer to store experience tuples."""

      def __init__(self, action_size, buffer_size, batch_size, seed):
          """Initialize a ReplayBuffer object.
          Params
          =====
              buffer_size (int): maximum size of buffer
              batch_size (int): size of each training batch
          """
          self.action_size = action_size
          self.memory = deque(maxlen=buffer_size)  # internal memory (deque)
          self.batch_size = batch_size
          self.experience = namedtuple("Experience", field_names=["state",
↪ "action", "reward", "next_state", "done"])
          self.seed = random.seed(seed)

      def add(self, state, action, reward, next_state, done):
          """Add a new experience to memory."""
          e = self.experience(state, action, reward, next_state, done)
          self.memory.append(e)

      def sample(self):
          """Randomly sample a batch of experiences from memory."""
          experiences = random.sample(self.memory, k=self.batch_size)

          states = torch.from_numpy(np.vstack([e.state for e in experiences if e
↪ is not None])).float().to(device)
          actions = torch.from_numpy(np.vstack([e.action for e in experiences if
↪ e is not None])).float().to(device)
          rewards = torch.from_numpy(np.vstack([e.reward for e in experiences if
↪ e is not None])).float().to(device)
          next_states = torch.from_numpy(np.vstack([e.next_state for e in
↪ experiences if e is not None])).float().to(device)
          dones = torch.from_numpy(np.vstack([e.done for e in experiences if e is
↪ not None])).astype(np.uint8).float().to(device)

          return (states, actions, rewards, next_states, dones)
```

```
def __len__(self):
    """Return the current size of internal memory."""
    return len(self.memory)
```

```
[13]: class OUNoise:
    """Ornstein-Uhlenbeck process."""

    def __init__(self, size, seed, mu=0., theta=0.15, sigma=0.2):
        """Initialize parameters and noise process."""
        self.mu = mu * np.ones(size)
        self.theta = theta
        self.sigma = sigma
        self.seed = random.seed(seed)
        self.size = size
        self.reset()

    def reset(self):
        """Reset the internal state (= noise) to mean (mu)."""
        self.state = copy.copy(self.mu)

    def sample(self):
        """Update internal state and return it as a noise sample."""
        x = self.state
        dx = self.theta * (self.mu - x) + self.sigma * np.random.
↪standard_normal(self.size)
        self.state = x + dx
        return self.state
```

```
[14]: class Agent():
    """Interacts with and learns from the environment."""

    def __init__(self, state_size, action_size, num_agents, random_seed):
        """Initialize an Agent object.

        Params
        =====
        state_size (int): dimension of each state
        action_size (int): dimension of each action
        random_seed (int): random seed
        """

        self.state_size = state_size
        self.action_size = action_size
        self.num_agents = num_agents
        self.seed = random.seed(random_seed)

        # Actor Network (w/ Target Network)
```



```

        self.actor_local = Actor(state_size, action_size, random_seed).
        ↪to(device)
        self.actor_target = Actor(state_size, action_size, random_seed).
        ↪to(device)
        self.actor_optimizer = optim.Adam(self.actor_local.parameters(),
        ↪lr=LR_ACTOR)

        # Critic Network (w/ Target Network)
        self.critic_local = Critic(state_size, action_size, random_seed).
        ↪to(device)
        self.critic_target = Critic(state_size, action_size, random_seed).
        ↪to(device)
        self.critic_optimizer = optim.Adam(self.critic_local.parameters(),
        ↪lr=LR_CRITIC, weight_decay=WEIGHT_DECAY)

        # Noise process
        self.noise = OUNoise((num_agents, action_size), random_seed)

        # Replay memory
        self.memory = ReplayBuffer(action_size, BUFFER_SIZE, BATCH_SIZE,
        ↪random_seed)

    def step(self, state, action, reward, next_state, done):
        """Save experience in replay memory, and use random sample from buffer
        ↪to learn."""
        # Save experience / reward
        for s, a, r, ns, d in zip(state, action, reward, next_state, done):
            self.memory.add(s, a, r, ns, d)

        # Learn, if enough samples are available in memory
        if len(self.memory) > BATCH_SIZE:
            experiences = self.memory.sample()
            self.learn(experiences, GAMMA)

    def act(self, state, add_noise=True):
        """Returns actions for given state as per current policy."""
        state = torch.from_numpy(state).float().to(device)
        self.actor_local.eval()
        with torch.no_grad():
            action = self.actor_local(state).cpu().data.numpy()
        self.actor_local.train()
        if add_noise:
            action += self.noise.sample()
        return np.clip(action, -1, 1)

    def reset(self):

```

```

self.noise.reset()

def learn(self, experiences, gamma):
    """Update policy and value parameters using given batch of experience_
    ↪ tuples.
    Q_targets = r + * critic_target(next_state, actor_target(next_state))
    where:
        actor_target(state) -> action
        critic_target(state, action) -> Q-value
    Params
    =====
    experiences (Tuple[torch.Tensor]): tuple of (s, a, r, s', done)_
    ↪ tuples
        gamma (float): discount factor
    """
    states, actions, rewards, next_states, dones = experiences

    # ----- update critic_
    ↪ ----- #
    # Get predicted next-state actions and Q values from target models
    actions_next = self.actor_target(next_states)
    Q_targets_next = self.critic_target(next_states, actions_next)
    # Compute Q targets for current states (y_i)
    Q_targets = rewards + (gamma * Q_targets_next * (1 - dones))
    # Compute critic loss
    Q_expected = self.critic_local(states, actions)
    critic_loss = F.mse_loss(Q_expected, Q_targets)
    # Minimize the loss
    self.critic_optimizer.zero_grad()
    critic_loss.backward()
    self.critic_optimizer.step()

    # ----- update actor_
    ↪ ----- #
    # Compute actor loss
    actions_pred = self.actor_local(states)
    actor_loss = -self.critic_local(states, actions_pred).mean()
    # Minimize the loss
    self.actor_optimizer.zero_grad()
    actor_loss.backward()
    self.actor_optimizer.step()

    # ----- update target networks_
    ↪ ----- #
    self.soft_update(self.critic_local, self.critic_target, TAU)
    self.soft_update(self.actor_local, self.actor_target, TAU)

```

```

def soft_update(self, local_model, target_model, tau):
    """Soft update model parameters.
    _target = *_local + (1 - )*_target
    Params
    =====
        local_model: PyTorch model (weights will be copied from)
        target_model: PyTorch model (weights will be copied to)
        tau (float): interpolation parameter
    """
    for target_param, local_param in zip(target_model.parameters(),
    ↪local_model.parameters()):
        target_param.data.copy_(tau*local_param.data + (1.
    ↪0-tau)*target_param.data)

```

1.0.7 Train

```

[15]: device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

if torch.cuda.is_available():
    print("CUDA Device:", torch.cuda.get_device_name(0))
    print("Memory Allocated:", round(torch.cuda.memory_allocated(0)/1024**3,1),
    ↪"GB")
else:
    print("Training on CPU")

```

CUDA Device: Tesla T4
Memory Allocated: 0.0 GB

```

[16]: # Initialize the agent
agent = Agent(state_size = state_size, action_size = action_size, num_agents =
    ↪num_agents, random_seed=1)

```

```

[17]: def train(n_episodes=300, print_every=50):
    scores_deque = deque(maxlen=print_every)
    scores = []
    for i_episode in range(1, n_episodes+1):
        # get initial environment and activate learning
        env_info = env.reset(train_mode=True)[brain_name]
        states = env_info.vector_observations
        agent.reset()
        score = 0.
        while True:
            actions = agent.act(states)
            env_info = env.step(actions)[brain_name]           # send all
            ↪actions to the environment
            next_states = env_info.vector_observations          # get next
            ↪state (for each agent)

```

```

        rewards = env_info.rewards                                # get
    ↪reward (for each agent)
        dones = env_info.local_done                               # see if
    ↪episode finished
        agent.step(states, actions, rewards, next_states, dones) # continue
    ↪to next step
        score += np.mean(env_info.rewards)                        # update
    ↪the score (for each agent)
        states = next_states                                      # roll
    ↪over states to next time step
        if np.any(dones):
            break

        scores_deque.append(score)
        scores.append(score)
        print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.
    ↪mean(scores_deque)), end="")
        if i_episode % 10 == 0:
            torch.save(agent.actor_local.state_dict(), 'checkpoint_actor.pth')
            torch.save(agent.critic_local.state_dict(), 'checkpoint_critic.pth')
        if i_episode % print_every == 0:
            print('\rEpisode {} \tAverage Score: {:.2f}'.format(i_episode, np.
    ↪mean(scores_deque)))

    return scores

scores = train()

```

```

Episode 50      Average Score: 20.44
Episode 100     Average Score: 36.64
Episode 150     Average Score: 36.12
Episode 200     Average Score: 36.39
Episode 250     Average Score: 36.28
Episode 300     Average Score: 35.70

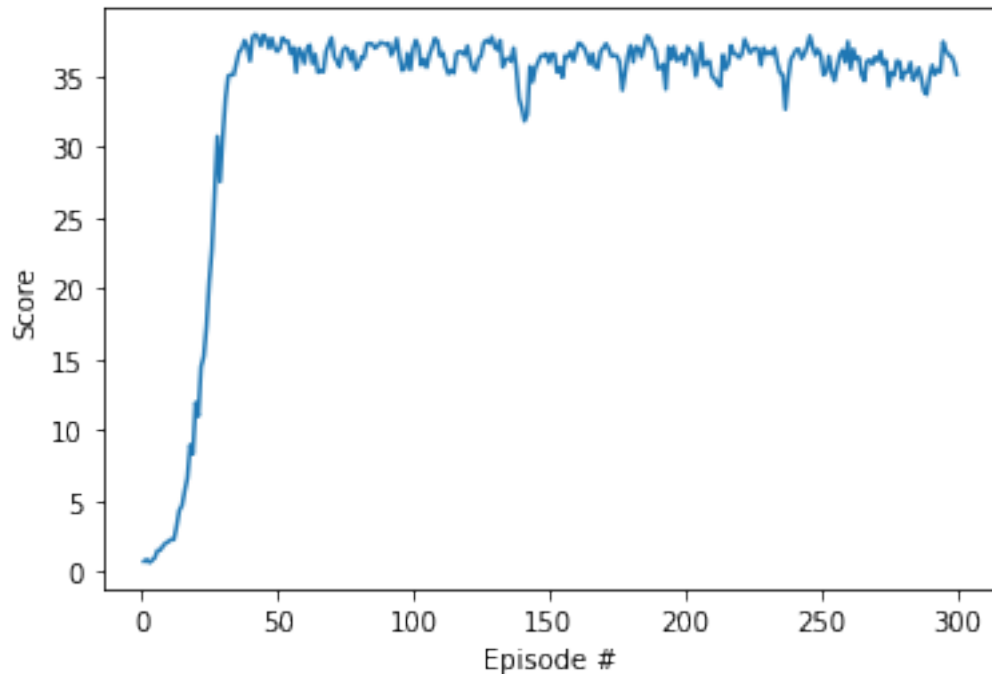
```

```

[18]: from matplotlib import pyplot as plt

fig = plt.figure()
ax = fig.add_subplot(111)
plt.plot(np.arange(1, len(scores)+1), scores)
plt.ylabel('Score')
plt.xlabel('Episode #')
plt.show()

```



1.0.8 Watch your Trained Agent play

```
[19]: # load models
actor_state_dict = torch.load('checkpoint_actor.pth')
print(actor_state_dict.keys())
agent.actor_local.load_state_dict(actor_state_dict)

critic_state_dict = torch.load('checkpoint_critic.pth')
print(critic_state_dict.keys())
agent.critic_local.load_state_dict(critic_state_dict)

odict_keys(['fc1.weight', 'fc1.bias', 'fc2.weight', 'fc2.bias', 'fc3.weight',
'fc3.bias'])
odict_keys(['fcs1.weight', 'fcs1.bias', 'fc2.weight', 'fc2.bias', 'fc3.weight',
'fc3.bias'])

[20]: env_info = env.reset(train_mode=False)[brain_name]      # reset the environment
      ↪
states = env_info.vector_observations                         # get the current state
      ↪ (for each agent)
scores = np.zeros(num_agents)                               # initialize the score
      ↪ (for each agent)

while True:
```

```

    actions = agent.act(states)
    env_info = env.step(actions)[brain_name]           # send all actions
    ↪to the environment
    next_states = env_info.vector_observations         # get next state
    ↪(for each agent)
    rewards = env_info.rewards                       # get reward (for
    ↪each agent)
    dones = env_info.local_done                      # see if episode
    ↪finished
    agent.step(states, actions, rewards, next_states, dones) # continue to next
    ↪step
    scores += env_info.rewards                        # update the score
    ↪(for each agent)
    states = next_states                             # roll over states
    ↪to next time step
    if np.any(dones):
        break

print('Total score (averaged over agents) this episode: {}'.format(np.
    ↪mean(scores)))

```

Total score (averaged over agents) this episode: 34.77849922263995

When finished, you can close the environment.

```
[21]: env.close()
```

1.0.9 Ideas for Future Work

- Spend more time on hyperparameter tuning
- Enhance Actor and Critic model architectures