

Report

July 2, 2021

1 Navigation

In this notebook, you will learn how to use the Unity ML-Agents environment.

1.1 Start the Environment

We begin by importing some 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
import torch
import torch.nn as nn
import torch.optim as optim
import torch.nn.functional as F
import random
from collections import namedtuple, deque
import matplotlib.pyplot as plt
%matplotlib inline
```

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/Banana.app"
- **Windows (x86)**: "path/to/Banana_Windows_x86/Banana.exe"
- **Windows (x86_64)**: "path/to/Banana_Windows_x86_64/Banana.exe"
- **Linux (x86)**: "path/to/Banana_Linux/Banana.x86"
- **Linux (x86_64)**: "path/to/Banana_Linux/Banana.x86_64"
- **Linux (x86, headless)**: "path/to/Banana_Linux_NoVis/Banana.x86"
- **Linux (x86_64, headless)**: "path/to/Banana_Linux_NoVis/Banana.x86_64"

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

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

```
[2]: env = UnityEnvironment(file_name="Banana.app")
```

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

Unity brain name: BananaBrain
    Number of Visual Observations (per agent): 0
    Vector Observation space type: continuous
    Vector Observation space size (per agent): 37
    Number of stacked Vector Observation: 1
    Vector Action space type: discrete
    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]
```

1.2 Examine the State and Action Spaces

The simulation contains a single agent that navigates a large environment. At each time step, it has four actions at its disposal: - 0 - walk forward - 1 - walk backward - 2 - turn left - 3 - turn right

The state space has 37 dimensions and contains the agent's velocity, along with ray-based perception of objects around agent's forward direction. A reward of +1 is provided for collecting a yellow banana, and a reward of -1 is provided for collecting a blue banana.

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 in the environment
print('Number of agents:', len(env_info.agents))

# number of actions
action_size = brain.vector_action_space_size
print('Number of actions:', action_size)

# examine the state space
state = env_info.vector_observations[0]
print('States look like:', state)
state_size = len(state)
print('States have length:', state_size)
```

```
Number of agents: 1
Number of actions: 4
States look like: [1.          0.          0.          0.          0.84408134 0.
 0.          1.          0.          0.0748472 0.          1.
 0.          0.          0.25755    1.          0.          0.
 0.          0.74177343 0.          1.          0.          0.
 0.25854847 0.          0.          1.          0.          0.09355672
 0.          1.          0.          0.          0.31969345 0.
 0.          ]
States have length: 37
```

1.3 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 (uniformly) 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
state = env_info.vector_observations[0] # get the current state
score = 0 # initialize the score
while True:
    action = np.random.randint(action_size) # select an action
    # send the action to the environment
    env_info = env.step(action)[brain_name]
    next_state = env_info.vector_observations[0] # get the next state
    reward = env_info.rewards[0] # get the reward
    done = env_info.local_done[0] # see if episode has finished
    score += reward # update the score
    # roll over the state to next time step
    state = next_state
    if done: # exit loop if episode finished
        break

print("Score: {}".format(score))
```

Score: 0.0

1.4 Define a neural network architecture that maps states to action values

Now it's your turn to train the 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]
```

1.4.1 Deep Q-Network (DQN)

Illustration of DQN Architecture

```
[6]: class QNetwork(nn.Module):
    """Actor (Policy) Model."""

    def __init__(self, state_size, action_size,
                 seed, fc1_units=64, fc2_units=64):
        """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(QNetwork, 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)

    def forward(self, state):
        """Build a network that maps state -> action values."""
        x = F.relu(self.fc1(state))
        x = F.relu(self.fc2(x))
        return self.fc3(x)
```

1.4.2 Fixed-size buffer to store experience tuples

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

```
device(type='cpu')
```

```

[8]: class ReplayBuffer:
    """Fixed-size buffer to store experience tuples."""
    def __init__(self, action_size, buffer_size, batch_size, seed):
        """Initialize a ReplayBuffer object.

        Params
        =====
            action_size (int): dimension of each action
            buffer_size (int): maximum size of buffer
            batch_size (int): size of each training batch
            seed (int): random seed
        """
        self.action_size = action_size
        self.memory = deque(maxlen=buffer_size)
        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]))\
            .long().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)

```

1.4.3 Define the Agent

The Agent Parameters

- **state_size** (int): Number of parameters in the environment state
- **action_size** (int): Number of actions
- **seed** (int): random seed
- **gamma** (float): discount factor

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

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

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

        # Q-Network
        self.qnetwork_local = QNetwork(state_size, action_size, seed).to(device)
        self.qnetwork_target = QNetwork(state_size, action_size, seed).to(device)
        self.optimizer = optim.Adam(self.qnetwork_local.parameters(), lr=LR)

        # Replay memory
        self.memory = ReplayBuffer(action_size, BUFFER_SIZE, BATCH_SIZE, seed)
        # Initialize time step (for updating every UPDATE_EVERY steps)
        self.t_step = 0

    def step(self, state, action, reward, next_state, done):
        # Save experience in replay memory
        self.memory.add(state, action, reward, next_state, done)

        # Learn every UPDATE_EVERY time steps.
        self.t_step = (self.t_step + 1) % UPDATE_EVERY
        if self.t_step == 0:
            # If enough samples are available in memory,
            # get random subset and learn
            if len(self.memory) > BATCH_SIZE:
                experiences = self.memory.sample()
                self.learn(experiences, GAMMA)
```

```

def act(self, state, eps=0.):
    """Returns actions for given state as per current policy.

    Params
    =====
        state (array_like): current state
        eps (float): epsilon, for epsilon-greedy action selection
    """

    state = torch.from_numpy(state).float().unsqueeze(0).to(device)
    self.qnetwork_local.eval()
    with torch.no_grad():
        action_values = self.qnetwork_local(state)
    self.qnetwork_local.train()

    # Epsilon-greedy action selection
    if random.random() > eps:
        return np.argmax(action_values.cpu().data.numpy())
    else:
        return random.choice(np.arange(self.action_size))

def learn(self, experiences, gamma):
    """Update value parameters using given batch of experience tuples.
    Params
    =====
        experiences (Tuple[torch.Variable]):
            tuple of (s, a, r, s', done) tuples
        gamma (float): discount factor
    """

    states, actions, rewards, next_states, dones = experiences
    # Get max predicted Q values (for next states) from target model
    Q_targets_next = self.qnetwork_target(next_states)\
        .detach().max(1)[0].unsqueeze(1)
    # Compute Q targets for current states
    Q_targets = rewards + (gamma * Q_targets_next * (1 - dones))
    # Get expected Q values from local model
    Q_expected = self.qnetwork_local(states).gather(1, actions)

    # Compute loss
    loss = F.mse_loss(Q_expected, Q_targets)
    # Minimize the loss
    self.optimizer.zero_grad()
    loss.backward()
    self.optimizer.step()

    # ----- update target network ----- #
    self.soft_update(self.qnetwork_local, self.qnetwork_target, TAU)

```



```

def soft_update(self, local_model, target_model, tau):
    """Soft update model parameters.
    theta_target = tau*theta_local + (1 - tau)*theta_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)

```

1.4.4 Hyperparameters

```

[10]: BUFFER_SIZE = int(1e5) # replay buffer size
      BATCH_SIZE = 64      # minibatch size
      GAMMA = 0.99        # discount factor
      TAU = 1e-3          # for soft update of target parameters
      LR = 5e-4           # learning rate
      UPDATE_EVERY = 4    # how often to update the network

```

1.5 Train the Agent with DQN

1.5.1 The Parameters

- **n_episodes** (int): Maximum number of training episodes
- **max_t** (int): Maximum number of timesteps per episode
- **eps_start** (float): starting value of epsilon, for epsilon-greedy action selection
- **eps_end** (float): minimum value of epsilon
- **eps_decay** (float): multiplicative factor (per episode) for decreasing epsilon

```

[11]: def dqn(n_episodes=2000, max_t=1000, eps_start=1.0,
             eps_end=0.01, eps_decay=0.995):
    """Deep Q-Learning.

    Params
    =====
        n_episodes (int) : maximum number of training episodes
        max_t (int)      : maximum number of timesteps per episode
        eps_start (float): starting value of epsilon,
                           for epsilon-greedy action selection
        eps_end (float)  : minimum value of epsilon
        eps_decay (float): multiplicative factor (per episode)
                           for decreasing epsilon

```

```

"""
scores = []                                # list containing scores from each episode
scores_window = deque(maxlen=100)          # last 100 scores
eps = eps_start                            # initialize epsilon

for i_episode in range(1, n_episodes+1):
    env_info = env.reset(train_mode=True)[brain_name] # reset the environment
    state = env_info.vector_observations[0]           # get the next state

    score = 0 # initialize the score

    for t in range(max_t):
        action = agent.act(state, eps)              # get action
        # send the action to the environment
        env_info = env.step(action)[brain_name]
        next_state = env_info.vector_observations[0] # get the next state
        reward = env_info.rewards[0]                # get the reward
        # see if episode has finished
        done = env_info.local_done[0]

        agent.step(state, action, reward, next_state, done)
        state = next_state
        score += reward
        if done or t > max_t:
            break

    scores_window.append(score)                    # save most recent score
    scores.append(score)                          # save most recent score
    eps = max(eps_end, eps_decay*eps) # decrease epsilon

    print('\rEpisode {} \tAverage Score: {:.2f}'.\
          format(i_episode, np.mean(scores_window)), end="")

    if i_episode % 100 == 0:
        print('\rEpisode {} \tAverage Score: {:.2f}'.\
              format(i_episode, np.mean(scores_window)))

    if np.mean(scores_window) >= 13.0:
        print('\nEnv. solved in {:d} episodes! \tAverage Score: {:.2f}'.\
              format(i_episode-100, np.mean(scores_window)))
        torch.save(agent.qnetwork_local.state_dict(), 'checkpoint.pth')
        break
return scores

```

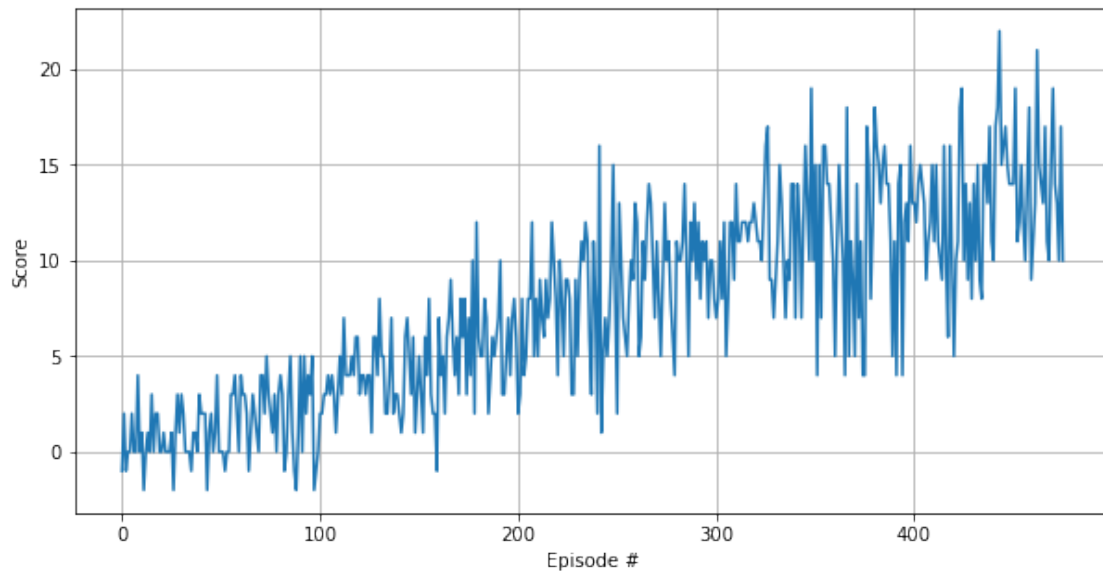
```
[12]: env_info = env.reset(train_mode=False)[brain_name] # reset the environment
```

```
[13]: state_size = len(env_info.vector_observations[0])  
      action_size = brain.vector_action_space_size # number of actions  
      agent = Agent(state_size=state_size, action_size=action_size, seed=0)
```

```
[14]: scores = dqn(n_episodes = 4000)
```

```
Episode 100      Average Score: 1.29  
Episode 200      Average Score: 4.58  
Episode 300      Average Score: 8.43  
Episode 400      Average Score: 11.29  
Episode 476      Average Score: 13.03  
Env. solved in 376 episodes!   Average Score: 13.03
```

```
[15]: # plot the scores  
fig = plt.figure(figsize=(10,5))  
plt.plot(np.arange(len(scores)), scores)  
plt.ylabel('Score')  
plt.xlabel('Episode #')  
plt.grid()  
plt.show()
```



When finished, you can close the environment.

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

1.6 Save model weights

```
[17]: torch.save(agent.qnetwork_local.state_dict(), 'model.pt')
```

1.7 Ideas for Future Work

Implement a double DQN (<https://arxiv.org/abs/1509.06461>), a dueling DQN (<https://arxiv.org/abs/1511.06581>), and/or prioritized experience replay (<https://arxiv.org/abs/1511.05952>)