


[Open in app](#)

Adrian Chow

[About](#)

Solving Continuous Control using Deep Reinforcement Learning (Policy-Based Methods)

 **Adrian Chow** Just now · 6 min read

Introduction

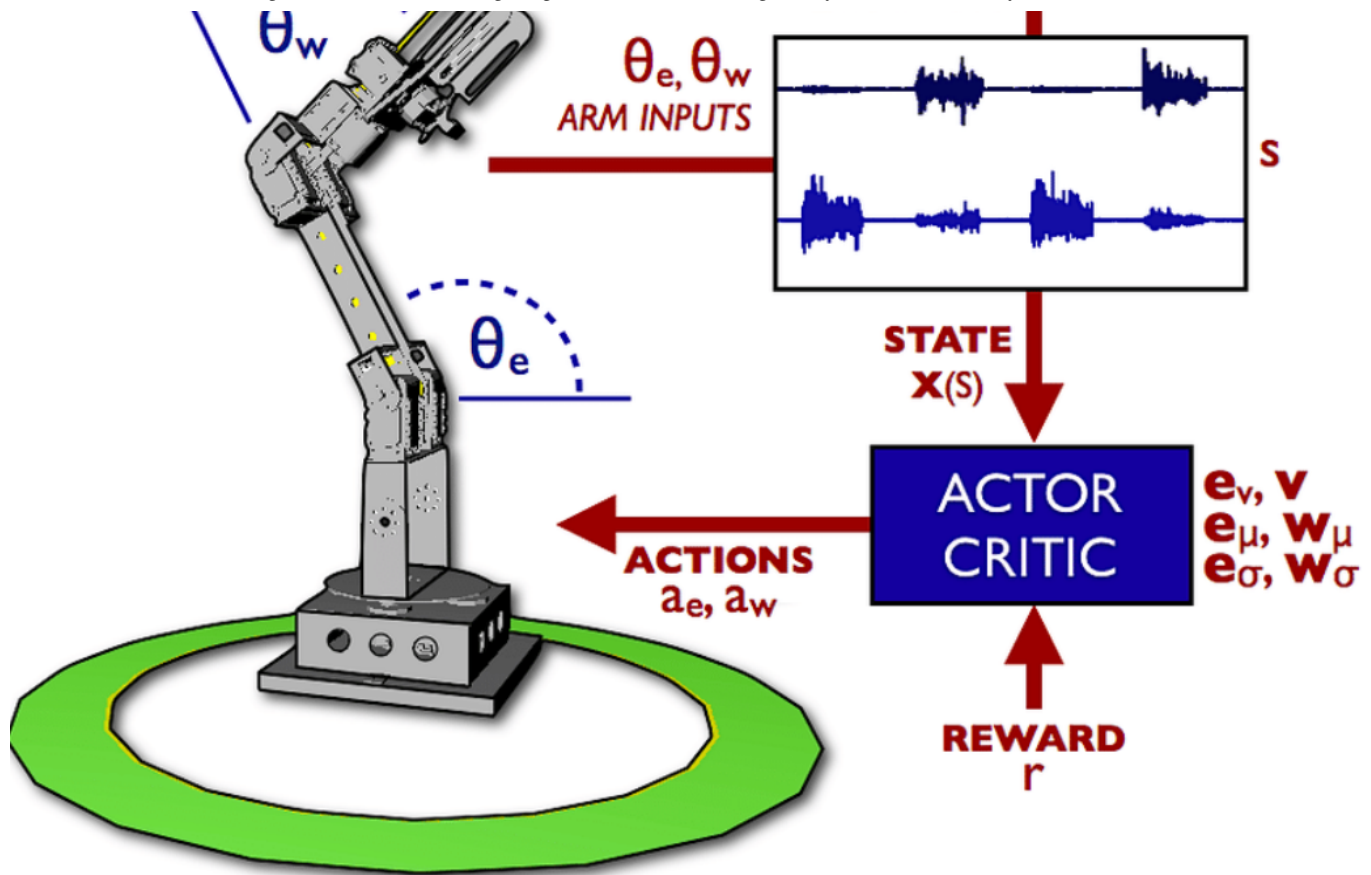
An alternative to classical control methods is deep reinforcement learning. Both are used to solve an optimization problem for dynamic systems that have a target behavior.

Classical control theory deals with the behavior of dynamical systems with inputs and how behavior can be tuned by using feedback. On the other hand, Deep RL's approach relies on an agent that is trained to have the policy which maximizes a measurable reward. The following article elaborates on a Deep RL agent's ability to solve a continuous control problem, namely Unity's Reacher.

Real-World Robotics

The application of this project is targeted for, while not limited to, robotic arms. Once the control aspect of the stack is solved, one can then generate a behavior to plan tasks.





Robotic Arm using Deep Reinforcement Learning (Source: ResearchGate)

Furthermore, it is possible to increase the productivity of manufacturing by sharing the trained policy. It has been shown that having multiple copies of the same agent sharing experience can accelerate learning, as I learned from solving the project!

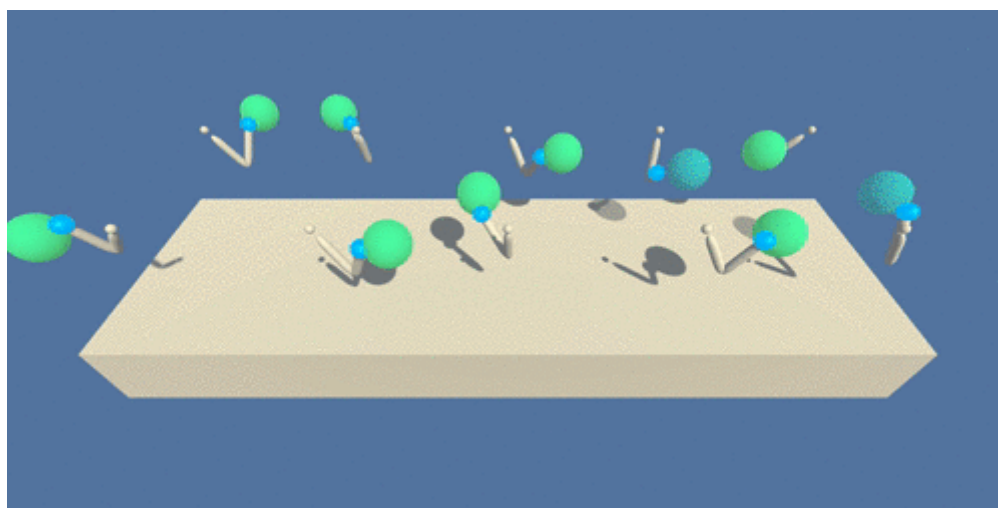


Multiple Robotic Arms using same Policy

The Environment

To solve the task of continuous control, I have used the Reacher environment on the [Unity ML-Agents GitHub page](#).

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.



Reacher Environment (Source: Unity ML-Agents)

RL Problem Specifications

- Goal of Agent: move to target locations
- Rewards : $+0.1$ every step in the goal location, $+0.0$ every step out of the goal location
- Action Space — Continuous, 4-Dimension Vector Space from $[-1, +1]$
- State Space — Continuous, 33-Dimension

- Solving Condition: Average score of +300.00 over 100 consecutive episodes.

Learning Algorithm

To solve the environment I have used an Actor-Critic Method which is a Deep Reinforcement Learning agent that utilizes two neural networks to estimate the policy (*actor*) and value function (*critic*). The policy structure is known as the *actor*, because it is used to select actions, and the estimated value function is known as the *critic* because it criticizes the actions made by the actor. Although I will go through the general algorithm, you can learn more about Actor-Critic Methods from [Chris Yoon's in-depth article about Actor-Critic Methods](#). His work was tremendously useful when exploring this topic.

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:

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$

end for

end for

DDPG Algorithm (Source: [Here](#))

The adaptation for the specific Actor-Critic method used is referred to as the Deep Deterministic Policy Gradient (DDPG). The original paper can be found [here](#). The idea

stems from the success of Deep-Q Learning and is modified to continuous action spaces.

The main takeaways are:

- Uses an Actor to choose the agent's actions deterministically.
- Uses a Critic to estimate the return value of the next state-action pair.
- DDPG is an Off-Policy Actor-Critic algorithm: Policy used to interact with the environment is different from the policy being learned.
- Uses soft-update to update target networks
- Uses some controlled action noise to explore action space

Initialization:

- Randomly initialize critic network $Q(s, a | \theta-Q)$ and actor $\mu(s | \theta-\mu)$ with weights $\theta-Q$ (critic) and $\theta-\mu$ (actor).
- Initialize target networks denoted by Q' and μ' with weights $\theta-Q' \leftarrow \theta-Q$, $\theta-\mu' \leftarrow \theta-\mu$
- Initialize replay buffer R . The replay buffer is a data structure to This data structure that stores nodes (named tuples) with the following data: ["state", "action", "reward", "next_state", "done"]

Here is the implementation for the Actor and Critic Neural Networks.

Note that:

- The Actor uses the **Hyperbolic Tangent Activation Function (tanh)** to limit the action space to $[-1, 1]$.
- Using a seed will result in the same weights for initialization between the target and local model.
- Critic factors in the action during the second fully-connected layer.

```

1  import numpy as np
2
3  import torch
4  import torch.nn as nn
5  import torch.nn.functional as F
6
7  def hidden_init(layer):
8      fan_in = layer.weight.data.size()[0]
9      lim = 1. / np.sqrt(fan_in)
10     return (-lim, lim)
11
12     class Actor(nn.Module):
13         """Actor (Policy) Model."""
14
15         def __init__(self, state_size, action_size, seed, fc1_units=128, fc2_units=128):
16             """Initialize parameters and build model.
17             Params
18             =====
19                 state_size (int): Dimension of each state
20                 action_size (int): Dimension of each action
21                 seed (int): Random seed
22                 fc1_units (int): Number of nodes in first hidden layer
23                 fc2_units (int): Number of nodes in second hidden layer
24             """
25             super(Actor, self).__init__()
26             self.seed = torch.manual_seed(seed)
27             self.fc1 = nn.Linear(state_size, fc1_units)
28             self.fc2 = nn.Linear(fc1_units, fc2_units)
29             self.fc3 = nn.Linear(fc2_units, action_size)
30             self.reset_parameters()
31
32         def reset_parameters(self):
33             self.fc1.weight.data.uniform_(*hidden_init(self.fc1))
34             self.fc2.weight.data.uniform_(*hidden_init(self.fc2))
35             self.fc3.weight.data.uniform_(-3e-3, 3e-3)
36
37         def forward(self, state):
38             """Build an actor (policy) network that maps states -> actions."""
39             x = F.relu(self.fc1(state))
40             x = F.relu(self.fc2(x))
41             return F.tanh(self.fc3(x))
42
43
44     class Critic(nn.Module):
45         """Critic (Value) Model """

```

```

45 critic (value) model.
46
47 def __init__(self, state_size, action_size, seed, fcs1_units=128, fc2_units=128):
48     """Initialize parameters and build model.
49     Params
50     =====
51     state_size (int): Dimension of each state
52     action_size (int): Dimension of each action
53     seed (int): Random seed
54     fcs1_units (int): Number of nodes in the first hidden layer
55     fc2_units (int): Number of nodes in the second hidden layer
56     """
57     super(Critic, self).__init__()
58     self.seed = torch.manual_seed(seed)
59     self.fcs1 = nn.Linear(state_size, fcs1_units)
60     self.fc2 = nn.Linear(fcs1_units+action_size, fc2_units)
61     self.fc3 = nn.Linear(fc2_units, 1)
62     self.reset_parameters()
63
64     def reset_parameters(self):
65         self.fcs1.weight.data.uniform_(*hidden_init(self.fcs1))
66         self.fc2.weight.data.uniform_(*hidden_init(self.fc2))
67         self.fc3.weight.data.uniform_(-3e-3, 3e-3)
68
69     def forward(self, state, action):
70         """Build a critic (value) network that maps (state, action) pairs -> Q-values."""
71         xs = F.relu(self.fcs1(state))
72         x = torch.cat((xs, action), dim=1)
73         x = F.relu(self.fc2(x))
74         return self.fc3(x)

```

actor_critic_implementation.py hosted with ❤ by GitHub

[view raw](#)

Here is the implementation of the Replay Buffer:

```

1  from collections import namedtuple, deque
2  import numpy as np
3  import random
4  import torch
5
6  field_names = ["state", "action", "reward", "next_state", "done"]
7
8  class ReplayBuffer:
9      """ Fixed-size buffer to store experience tuples"""

```



```
10
11 def __init__(self, action_size, buffer_size, batch_size, seed, device):
12     """Initialize a ReplayBuffer object. """
13     self.action_size = action_size
14     self.buffer_size = buffer_size # size of replay buffer
15     self.batch_size = batch_size # how many mem tuples to sample at a time
16     self.seed = random.seed(seed)
17     self.device = device
18
19     # Define Named Tuple - field_names=["state", "action", "reward", "next_state", "
20     self.experience = namedtuple("Experience", field_names=field_names)
21
22     # Data structure to hold the memories
23     self.memory = deque(maxlen=buffer_size)
24
25 def add(self, state, action, reward, next_state, done):
26     """Add a new experience to memory."""
27     e = self.experience(state, action, reward, next_state, done)
28     self.memory.append(e)
29
30
31 def sample(self):
32     """ Randomly sample a batch of experiences """
33
34     # Sample an experience with length k from list of memories
35     experiences = random.sample(self.memory, k=self.batch_size)
36
37     # For each item in the tuple, stack vertically and convert to GPU torch tensor
38     states = np.vstack([e.state for e in experiences if e is not None])
39     states = torch.from_numpy(states).float().to(self.device) # (float)
40
41     actions = np.vstack([e.action for e in experiences if e is not None])
42     actions = torch.from_numpy(actions).float().to(self.device) # (float)
43
44     rewards = np.vstack([e.reward for e in experiences if e is not None])
45     rewards = torch.from_numpy(rewards).float().to(self.device) # (float)
46
47     next_states = np.vstack([e.next_state for e in experiences if e is not None])
48     next_states = torch.from_numpy(next_states).float().to(self.device) # float
49
50     dones = np.vstack([e.done for e in experiences if e is not None]).astype(np.uint
51     dones = torch.from_numpy(dones).float().to(self.device)
52
53     return (states, actions, rewards, next_states, dones)
54
```



```

55
56     def __len__(self):
57         """Return the current size of internal memory."""
58         return len(self.memory)

```

replay_buffer.py hosted with ❤ by GitHub

[view raw](#)

Putting it together:

```

1  import numpy as np
2  import random
3  import copy
4  from collections import namedtuple, deque
5
6  from networks import Actor, Critic
7
8  import torch
9  import torch.nn.functional as F
10 import torch.optim as optim
11
12 BUFFER_SIZE = int(1e5) # replay buffer size
13 BATCH_SIZE = 128      # minibatch size
14 GAMMA = 0.99          # discount factor
15 TAU = 1e-3            # for soft update of target parameters
16 LR_ACTOR = 1e-3       # learning rate of the actor
17 LR_CRITIC = 1e-3      # learning rate of the critic
18 WEIGHT_DECAY = 0      # L2 weight decay
19
20 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
21
22 class DDPG_Agent():
23     """Interacts with and learns from the environment."""
24
25     def __init__(self, state_size, action_size, n_agents, random_seed):
26         """Initialize an Agent object.
27
28         Params
29         =====
30             state_size (int): dimension of each state
31             action_size (int): dimension of each action
32             random_seed (int): random seed
33         """
34         self.state_size = state_size
35         self.action_size = action_size

```

```

35 self.action_size = action_size
36 self.seed = random.seed(random_seed)
37
38 # Actor Network (w/ Target Network)
39 self.actor_local = Actor(state_size, action_size, random_seed).to(device)
40 self.actor_target = Actor(state_size, action_size, random_seed).to(device)
41 self.actor_optimizer = optim.Adam(self.actor_local.parameters(), lr=LR_ACTOR)
42
43 # Critic Network (w/ Target Network)
44 self.critic_local = Critic(state_size, action_size, random_seed).to(device)
45 self.critic_target = Critic(state_size, action_size, random_seed).to(device)
46 self.critic_optimizer = optim.Adam(self.critic_local.parameters(), lr=LR_CRITIC,
47
48 # Noise process
49 self.noise = OUNoise((n_agents, action_size), random_seed)
50
51 # Replay memory
52 self.memory = ReplayBuffer(action_size, BUFFER_SIZE, BATCH_SIZE, random_seed)

```

DDPG_init.py hosted with ❤ by GitHub

[view raw](#)

Training Loop

```

1  def follow_goal_ddpg(agent,
2      env,
3      brain_name,
4      n_agents,
5      n_episodes=1500,
6      max_t=3000,
7      print_every=10,
8      win_condition=30.0):
9
10     """
11
12     Continious Control using DDPG.
13
14     Params
15     =====
16     n_episodes (int): maximum number of training episodes
17     max_t (int): maximum number of timesteps per episode
18
19     """

```

```

18     print_every (int): how many episodes before printing scores
19     n_agents (int): how many arms are in the environment
20
21     """
22
23     scores = []
24     scores_mean = []
25     scores_window = deque(maxlen=100) # Score last 100 scores
26
27     for i_episode in tqdm(range(1, n_episodes+1)):
28
29         env_info = env.reset(train_mode=True)[brain_name] # reset the environment
30         states = env_info.vector_observations                # get the current state (for each agent)
31         score = np.zeros(n_agents)                         # initialize the score (for each agent)
32
33         for t in range(max_t):
34             actions = agent.act(states)                     # consult agent for actions
35             env_info = env.step(actions)[brain_name]        # take a step in the env
36             next_states = env_info.vector_observations       # get next state (for each agent)
37             rewards = env_info.rewards                      # get reward (for each agent)
38             dones = env_info.local_done                    # see if episode finished
39
40             # take a learning step
41             agent.step(states, actions, rewards, next_states, dones)
42
43             score += env_info.rewards                      # update the score (for each agent)
44             states = next_states                          # roll over states to next time step
45             if np.any(dones):                             # exit loop if episode finished
46                 break
47
48         scores.append(np.mean(score))
49         scores_window.append(np.mean(score))
50         scores_mean.append(np.mean(scores_window))
51
52     # Print on print_every condition
53     if i_episode % print_every == 0:
54         print('\rEpisode {} \tScore: {:.2f} \tAverage Score: {:.2f}'.format(i_episode,
55                                     np.mean(scores_window), np.mean(scores_mean)))
56
57     # Winning condition + save model parameters
58     if np.mean(scores_window) >= win_condition:
59         print('\nEnvironment solved in {:d} episodes! \tScore: {:.2f} \tAverage Score: {:.2f}'.format(
60             i_episode, np.mean(scores_window), np.mean(scores_mean)))
61         torch.save(agent.actor_local.state_dict(), 'checkpoint_actor.pth')
62         torch.save(agent.critic_local.state_dict(), 'checkpoint_critic.pth')
63         break

```

```

63     execution_info = {'last_score': scores.pop(),
64                       'solved_in': i_episode,
65                       'last_100_avg': np.mean(scores_window),
66                       'save_file': 'DDPG/...'}
67
68     return scores, scores_mean, execution_info

```

Training_loop.py hosted with ❤ by GitHub

[view raw](#)

For episode $e \leftarrow 1$ to M :

- Initialize a random process N for action exploration
- Receive initial observation state, s_1

```

for i_episode in range(1, n_episodes + 1):

    env_info = env.reset(train_mode=True)[brain_name]
    states = env_info.vector_observations
    score = np.zeros(n_agents)

```

Here we are able to sample our initial environment state as well as set up a list to store our training scores.

For step $t \leftarrow 1$ to T :

```

for t in range(max_t):
    actions = agent.act(states)
    env_info = env.step(actions)[brain_name]
    next_states = env_info.vector_observations
    rewards = env_info.rewards
    dones = env_info.local_done
    agent.step(states, actions, rewards, next_states, dones)
    score += env_info.rewards
    states = next_states
    if np.any(dones): # Check if there are any done agents
        break

```

- Select action $A = \mu(st|\theta - \mu) + \text{Noise}$ according to the current policy and exploration noise

```

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

```

get_action.py hosted with ❤ by GitHub

[view raw](#)

Here we use the local Actor model to sample an action space with a bit of added noise exploration. For the noise added to the actions, DDPGs often use the Ornstein–Uhlenbeck process to generate temporally correlated exploration for exploration efficiency in physical control problems.

```

1  class OUNoise:
2      """Ornstein-Uhlenbeck process."""
3
4      def __init__(self, size, seed, mu=0., theta=0.15, sigma=0.2):
5          """Initialize parameters and noise process."""
6          self.size = size
7          self.mu = mu * np.ones(size)
8          self.theta = theta
9          self.sigma = sigma
10         self.seed = random.seed(seed)
11         self.reset()
12
13     def reset(self):
14         """Reset the internal state (= noise) to mean (mu)."""
15         self.state = copy.copy(self.mu)
16
17     def sample(self):
18         """Update internal state and return it as a noise sample."""
19         x = self.state
20         dx = self.theta * (self.mu - x) + self.sigma * np.random.standard_normal(self.size)
21         self.state = x + dx
22         return self.state

```

Ornstein_Uhlenbeck_Process.py hosted with ❤ by GitHub

[view raw](#)

- Execute action at and observe reward $r(t)$ and observe new state $s(t+1)$, where t is the current timestep

```
env_info = env.step(actions)[brain_name]
next_states = env_info.vector_observations
rewards = env_info.rewards
dones = env_info.local_done
```

- Store transition $s(t), a(t), r(t), s(t+1)$ in Replay Buffer, R

```
agent.step(states, actions, rewards, next_states, dones)
```

```
1  def step(self, states, actions, rewards, next_states, dones):
2      """Save experience in replay memory, and use random sample from buffer to learn."""
3      # Save experience / reward
4
5      #for i in range(len(states)):
6      #    self.memory.add(states[i, :], actions[i, :], rewards[i], next_states[i, :], done[i])
7      for state, action, reward, next_state, done in zip(states, actions, rewards, next_states, dones):
8          self.memory.add(state, action, reward, next_state, done)
9
10     # Learn, if enough samples are available in memory
11     if len(self.memory) > BATCH_SIZE:
12         experiences = self.memory.sample()
13         self.learn(experiences, GAMMA)
```

step.py hosted with ❤ by GitHub

[view raw](#)

- The remaining parts are the algorithm are mathematical learning steps to update the policy (weights of the local and target networks).

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

~~Update the target networks.~~

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

$$\theta^{\mu'} \leftarrow \tau \theta^{\mu} + (1 - \tau) \theta^{\mu'}$$

end for
end for

Learning Steps of DDPG Algorithm

The agent takes a step and immediately samples randomly from the replay buffer to learn a bit more about the environment it is currently in. We can split the update into two parts: Actor and Critic update. (Note: Actor loss is -ve due to the fact we are trying to maximize the return values passed by the Critic.)

DDPG only updates the local networks (the networks that interact with the environment) and performs soft updates on the target networks. As you can see below, soft updates only factor in a small amount of the local network weights (defined by Tau).

```

1  def learn(self, experiences, gamma):
2      """Update policy and value parameters using given batch of experience tuples.
3      Q_targets = r + γ * critic_target(next_state, actor_target(next_state))
4      where:
5          actor_target(state) -> action
6          critic_target(state, action) -> Q-value
7      Params
8      =====
9          experiences (Tuple[torch.Tensor]): tuple of (s, a, r, s', done) tuples
10         gamma (float): discount factor
11     """
12     states, actions, rewards, next_states, dones = experiences
13
14     # ----- update critic ----- #
15     # Get predicted next-state actions and Q values from target models
16     actions_next = self.actor_target(next_states)
17     Q_targets_next = self.critic_target(next_states, actions_next)
18     # Compute Q targets for current states (y_i)
19     Q_targets = rewards + (gamma * Q_targets_next * (1 - dones))
20     # Compute critic loss
21     Q_expected = self.critic_local(states, actions)
22     critic_loss = F.mse_loss(Q_expected, Q_targets)
23     # Minimize the loss
24     self.critic_optimizer.zero_grad()
25     critic_loss.backward()
26     self.critic_optimizer.step()

```



```

27
28 # ----- update actor ----- #
29 # Compute actor loss
30 actions_pred = self.actor_local(states)
31 actor_loss = -self.critic_local(states, actions_pred).mean()
32 # Minimize the loss
33 self.actor_optimizer.zero_grad()
34 actor_loss.backward()
35 self.actor_optimizer.step()
36
37 # ----- update target networks ----- #
38 self.soft_update(self.critic_local, self.critic_target, TAU)
39 self.soft_update(self.actor_local, self.actor_target, TAU)
40
41 def soft_update(self, local_model, target_model, tau):
42     """Soft update model parameters.
43      $\theta_{\text{target}} = \tau \theta_{\text{local}} + (1 - \tau) \theta_{\text{target}}$ 
44     Params
45     =====
46         local_model: PyTorch model (weights will be copied from)
47         target_model: PyTorch model (weights will be copied to)
48         tau (float): interpolation parameter
49     """
50     for target_param, local_param in zip(target_model.parameters(), local_model.parameters()):
51         target_param.data.copy_(tau*local_param.data + (1.0-tau)*target_param.data)

```

learn.py hosted with ❤ by GitHub

[view raw](#)

Yeah, so that's it! The model is able to learn over time and develop the ability to follow a goal location in a continuously controlled state. The full repository can be viewed [here](#).

Plot of Results

Training Parameters:

- Max Episodes: 1500
- Max Time Steps: 3000
- Buffer Size: 10000
- Batch Size: 128