

<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div><div>Keras</div></div>
About Keras
Getting started
Developer guides
Keras API reference
Code examples
Why choose Keras?
Community & governance
Contributing to Keras

» [Code examples](#) / [Reinforcement learning](#) / Deep Deterministic Policy Gradient (DDPG)

Deep Deterministic Policy Gradient (DDPG)

Author: [amifunny](#)
Date created: 2020/06/04
Last modified: 2020/09/21
Description: Implementing DDPG algorithm on the Inverted Pendulum Problem.

 [View in Colab](#) ·  [GitHub source](#)

Introduction

Deep Deterministic Policy Gradient (DDPG) is a model-free off-policy algorithm for learning continous actions.

It combines ideas from DPG (Deterministic Policy Gradient) and DQN (Deep Q-Network). It uses Experience Replay and slow-learning target networks from DQN, and it is based on DPG, which can operate over continuous action spaces.

This tutorial closely follow this paper - [Continuous control with deep reinforcement learning](#)

Problem

We are trying to solve the classic **Inverted Pendulum** control problem. In this setting, we can take only two actions: swing left or swing right.

What make this problem challenging for Q-Learning Algorithms is that actions are **continuous** instead of being **discrete**. That is, instead of using two discrete actions like **-1** or **+1**, we have to select from infinite actions ranging from **-2** to **+2**.

Quick theory

Just like the Actor-Critic method, we have two networks:

- 1. Actor - It proposes an action given a state.
- 2. Critic - It predicts if the action is good (positive value) or bad (negative value) given a state and an action.

DDPG uses two more techniques not present in the original DQN:

First, it uses two Target networks.

Why? Because it add stability to training. In short, we are learning from estimated targets and Target networks are updated slowly, hence keeping our estimated targets stable.

Conceptually, this is like saying, "I have an idea of how to play this well, I'm going to try it out for a bit until I find something better", as opposed to saying "I'm going to re-learn how to play this entire game after every move". See this [StackOverflow answer](#).

Second, it uses Experience Replay.

We store list of tuples (**state**, **action**, **reward**, **next_state**), and instead of learning only from recent experience, we learn from sampling all of our experience accumulated so far.

Now, let's see how is it implemented.

```
import gym
import tensorflow as tf
from tensorflow.keras import layers
import numpy as np
import matplotlib.pyplot as plt
```

We use [OpenAI Gym](#) to create the environment. We will use the `upper_bound` parameter to scale our actions later.

```
problem = "Pendulum-v0"
env = gym.make(problem)

num_states = env.observation_space.shape[0]
print("Size of State Space -> {}".format(num_states))
num_actions = env.action_space.shape[0]
print("Size of Action Space -> {}".format(num_actions))

upper_bound = env.action_space.high[0]
lower_bound = env.action_space.low[0]

print("Max Value of Action -> {}".format(upper_bound))
print("Min Value of Action -> {}".format(lower_bound))
```

```
Size of State Space -> 3
Size of Action Space -> 1
Max Value of Action -> 2.0
Min Value of Action -> -2.0
```

To implement better exploration by the Actor network, we use noisy perturbations, specifically an **Ornstein-Uhlenbeck process** for generating noise, as described in the paper. It samples noise from a correlated normal distribution.

```
class OUActionNoise:
    def __init__(self, mean, std_deviation, theta=0.15, dt=1e-2, x_initial=None):
        self.theta = theta
        self.mean = mean
        self.std_dev = std_deviation
        self.dt = dt
        self.x_initial = x_initial
        self.reset()

    def __call__(self):
        # Formula taken from https://www.wikipedia.org/wiki/Ornstein-Uhlenbeck_process.
        x = (
            self.x_prev
            + self.theta * (self.mean - self.x_prev) * self.dt
            + self.std_dev * np.sqrt(self.dt) * np.random.normal(size=self.mean.shape)
        )
        # Store x into x_prev
        # Makes next noise dependent on current one
        self.x_prev = x
        return x

    def reset(self):
        if self.x_initial is not None:
            self.x_prev = self.x_initial
        else:
            self.x_prev = np.zeros_like(self.mean)
```

The `Buffer` class implements Experience Replay.

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 $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'}))$
 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

Critic loss - Mean Squared Error of $y - Q(s, a)$ where y is the expected return as seen by the Target network, and $Q(s, a)$ is action value predicted by the Critic network. y is a moving target that the critic model tries to achieve; we make this target stable by updating the Target model slowly.

Actor loss - This is computed using the mean of the value given by the Critic network for the actions taken by the Actor network. We seek to maximize this quantity.

Hence we update the Actor network so that it produces actions that get the maximum predicted value as seen by the Critic, for a given state.

```

class Buffer:
    def __init__(self, buffer_capacity=100000, batch_size=64):
        # Number of "experiences" to store at max
        self.buffer_capacity = buffer_capacity
        # Num of tuples to train on.
        self.batch_size = batch_size

        # Its tells us num of times record() was called.
        self.buffer_counter = 0

        # Instead of list of tuples as the exp.replay concept go
        # We use different np.arrays for each tuple element
        self.state_buffer = np.zeros((self.buffer_capacity, num_states))
        self.action_buffer = np.zeros((self.buffer_capacity, num_actions))
        self.reward_buffer = np.zeros((self.buffer_capacity, 1))
        self.next_state_buffer = np.zeros((self.buffer_capacity, num_states))

    # Takes (s,a,r,s') observation tuple as input
    def record(self, obs_tuple):
        # Set index to zero if buffer_capacity is exceeded,
        # replacing old records
        index = self.buffer_counter % self.buffer_capacity

        self.state_buffer[index] = obs_tuple[0]
        self.action_buffer[index] = obs_tuple[1]
        self.reward_buffer[index] = obs_tuple[2]
        self.next_state_buffer[index] = obs_tuple[3]

        self.buffer_counter += 1

    # Eager execution is turned on by default in TensorFlow 2. Decorating with tf.function
    allows
    # TensorFlow to build a static graph out of the logic and computations in our function.
    # This provides a large speed up for blocks of code that contain many small TensorFlow
    operations such as this one.
    @tf.function
    def update(
        self, state_batch, action_batch, reward_batch, next_state_batch,
    ):
        # Training and updating Actor & Critic networks.
        # See Pseudo Code.
        with tf.GradientTape() as tape:
            target_actions = target_actor(next_state_batch, training=True)
            y = reward_batch + gamma * target_critic(
                [next_state_batch, target_actions], training=True
            )
            critic_value = critic_model([state_batch, action_batch], training=True)
            critic_loss = tf.math.reduce_mean(tf.math.square(y - critic_value))

            critic_grad = tape.gradient(critic_loss, critic_model.trainable_variables)
            critic_optimizer.apply_gradients(
                zip(critic_grad, critic_model.trainable_variables)
            )

            with tf.GradientTape() as tape:
                actions = actor_model(state_batch, training=True)
                critic_value = critic_model([state_batch, actions], training=True)
                # Used `-value` as we want to maximize the value given
                # by the critic for our actions
                actor_loss = -tf.math.reduce_mean(critic_value)

            actor_grad = tape.gradient(actor_loss, actor_model.trainable_variables)
            actor_optimizer.apply_gradients(
                zip(actor_grad, actor_model.trainable_variables)
            )

    # We compute the loss and update parameters
    def learn(self):
        # Get sampling range
        record_range = min(self.buffer_counter, self.buffer_capacity)
        # Randomly sample indices
        batch_indices = np.random.choice(record_range, self.batch_size)

        # Convert to tensors
        state_batch = tf.convert_to_tensor(self.state_buffer[batch_indices])
        action_batch = tf.convert_to_tensor(self.action_buffer[batch_indices])
        reward_batch = tf.convert_to_tensor(self.reward_buffer[batch_indices])
        reward_batch = tf.cast(reward_batch, dtype=tf.float32)
        next_state_batch = tf.convert_to_tensor(self.next_state_buffer[batch_indices])

```

```

        self.update(state_batch, action_batch, reward_batch, next_state_batch)

# This update target parameters slowly
# Based on rate `tau`, which is much less than one.
@tf.function
def update_target(target_weights, weights, tau):
    for (a, b) in zip(target_weights, weights):
        a.assign(b * tau + a * (1 - tau))

```

Here we define the Actor and Critic networks. These are basic Dense models with **ReLU** activation.

Note: We need the initialization for last layer of the Actor to be between **-0.003** and **0.003** as this prevents us from getting **1** or **-1** output values in the initial stages, which would squash our gradients to zero, as we use the **tanh** activation.

```

def get_actor():
    # Initialize weights between -3e-3 and 3-e3
    last_init = tf.random_uniform_initializer(minval=-0.003, maxval=0.003)

    inputs = layers.Input(shape=(num_states,))
    out = layers.Dense(256, activation="relu")(inputs)
    out = layers.Dense(256, activation="relu")(out)
    outputs = layers.Dense(1, activation="tanh", kernel_initializer=last_init)(out)

    # Our upper bound is 2.0 for Pendulum.
    outputs = outputs * upper_bound
    model = tf.keras.Model(inputs, outputs)
    return model

def get_critic():
    # State as input
    state_input = layers.Input(shape=(num_states,))
    state_out = layers.Dense(16, activation="relu")(state_input)
    state_out = layers.Dense(32, activation="relu")(state_out)

    # Action as input
    action_input = layers.Input(shape=(num_actions,))
    action_out = layers.Dense(32, activation="relu")(action_input)

    # Both are passed through separate layer before concatenating
    concat = layers.Concatenate()([state_out, action_out])

    out = layers.Dense(256, activation="relu")(concat)
    out = layers.Dense(256, activation="relu")(out)
    outputs = layers.Dense(1)(out)

    # Outputs single value for give state-action
    model = tf.keras.Model([state_input, action_input], outputs)

    return model

```

policy() returns an action sampled from our Actor network plus some noise for exploration.

```

def policy(state, noise_object):
    sampled_actions = tf.squeeze(actor_model(state))
    noise = noise_object()
    # Adding noise to action
    sampled_actions = sampled_actions.numpy() + noise

    # We make sure action is within bounds
    legal_action = np.clip(sampled_actions, lower_bound, upper_bound)

    return [np.squeeze(legal_action)]

```

Training hyperparameters

```
std_dev = 0.2
ou_noise = OUActionNoise(mean=np.zeros(1), std_deviation=float(std_dev) * np.ones(1))

actor_model = get_actor()
critic_model = get_critic()

target_actor = get_actor()
target_critic = get_critic()

# Making the weights equal initially
target_actor.set_weights(actor_model.get_weights())
target_critic.set_weights(critic_model.get_weights())

# Learning rate for actor-critic models
critic_lr = 0.002
actor_lr = 0.001

critic_optimizer = tf.keras.optimizers.Adam(critic_lr)
actor_optimizer = tf.keras.optimizers.Adam(actor_lr)

total_episodes = 100
# Discount factor for future rewards
gamma = 0.99
# Used to update target networks
tau = 0.005

buffer = Buffer(50000, 64)
```

Now we implement our main training loop, and iterate over episodes. We sample actions using `policy()` and train with `learn()` at each time step, along with updating the Target networks at a rate `tau`.

```

# To store reward history of each episode
ep_reward_list = []
# To store average reward history of last few episodes
avg_reward_list = []

# Takes about 4 min to train
for ep in range(total_episodes):

    prev_state = env.reset()
    episodic_reward = 0

    while True:
        # Uncomment this to see the Actor in action
        # But not in a python notebook.
        # env.render()

        tf_prev_state = tf.expand_dims(tf.convert_to_tensor(prev_state), 0)

        action = policy(tf_prev_state, ou_noise)
        # Recieve state and reward from environment.
        state, reward, done, info = env.step(action)

        buffer.record((prev_state, action, reward, state))
        episodic_reward += reward

        buffer.learn()
        update_target(target_actor.variables, actor_model.variables, tau)
        update_target(target_critic.variables, critic_model.variables, tau)

        # End this episode when `done` is True
        if done:
            break

        prev_state = state

    ep_reward_list.append(episodic_reward)

    # Mean of last 40 episodes
    avg_reward = np.mean(ep_reward_list[-40:])
    print("Episode * {} * Avg Reward is ==> {}".format(ep, avg_reward))
    avg_reward_list.append(avg_reward)

# Plotting graph
# Episodes versus Avg. Rewards
plt.plot(avg_reward_list)
plt.xlabel("Episode")
plt.ylabel("Avg. Epsiodic Reward")
plt.show()

```



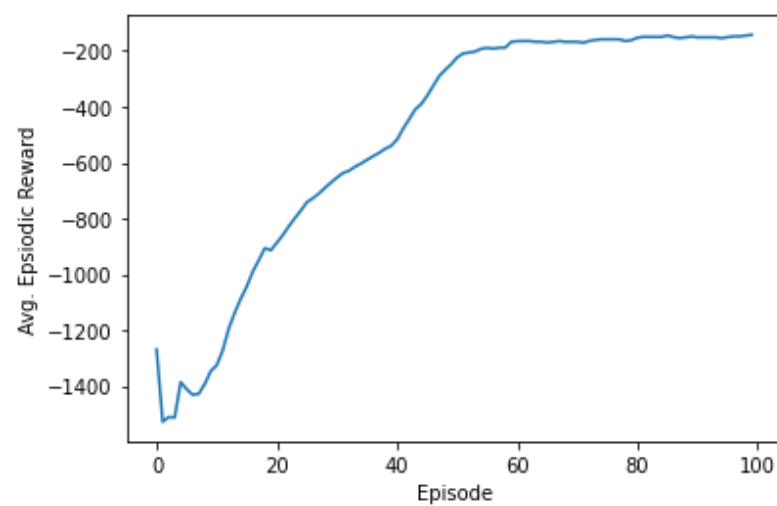
```
Episode * 0 * Avg Reward is ==> -1269.3278950595395
Episode * 1 * Avg Reward is ==> -1528.3008939716287
Episode * 2 * Avg Reward is ==> -1511.1737868279706
Episode * 3 * Avg Reward is ==> -1512.8568141261057
Episode * 4 * Avg Reward is ==> -1386.054573343386
Episode * 5 * Avg Reward is ==> -1411.4818856846339
Episode * 6 * Avg Reward is ==> -1431.6790621961388
Episode * 7 * Avg Reward is ==> -1427.9515009474867
Episode * 8 * Avg Reward is ==> -1392.9313930075857
Episode * 9 * Avg Reward is ==> -1346.6839043846012
Episode * 10 * Avg Reward is ==> -1325.5818224096574
Episode * 11 * Avg Reward is ==> -1271.778361283553
Episode * 12 * Avg Reward is ==> -1194.0784354001732
Episode * 13 * Avg Reward is ==> -1137.1096928093427
Episode * 14 * Avg Reward is ==> -1087.2426176918214
Episode * 15 * Avg Reward is ==> -1043.5265287176114
Episode * 16 * Avg Reward is ==> -990.0857409180443
Episode * 17 * Avg Reward is ==> -949.0661362879348
Episode * 18 * Avg Reward is ==> -906.1744575963231
Episode * 19 * Avg Reward is ==> -914.0098344966382
Episode * 20 * Avg Reward is ==> -886.8905055354011
Episode * 21 * Avg Reward is ==> -859.3416389004793
Episode * 22 * Avg Reward is ==> -827.5405203616622
Episode * 23 * Avg Reward is ==> -798.3875178404127
Episode * 24 * Avg Reward is ==> -771.289491103158
Episode * 25 * Avg Reward is ==> -741.6622445749622
Episode * 26 * Avg Reward is ==> -727.7080867854874
Episode * 27 * Avg Reward is ==> -710.485046117201
Episode * 28 * Avg Reward is ==> -690.3850022530833
Episode * 29 * Avg Reward is ==> -671.3205042911178
Episode * 30 * Avg Reward is ==> -653.4475135842247
Episode * 31 * Avg Reward is ==> -637.0057392119055
Episode * 32 * Avg Reward is ==> -629.2474166794424
Episode * 33 * Avg Reward is ==> -614.4655398230501
Episode * 34 * Avg Reward is ==> -603.3854873345723
Episode * 35 * Avg Reward is ==> -589.86534490467
Episode * 36 * Avg Reward is ==> -577.1806480684269
Episode * 37 * Avg Reward is ==> -565.1365286280546
Episode * 38 * Avg Reward is ==> -550.6647028563134
Episode * 39 * Avg Reward is ==> -540.0095147571197
Episode * 40 * Avg Reward is ==> -517.3861294233157
Episode * 41 * Avg Reward is ==> -478.705352005952
Episode * 42 * Avg Reward is ==> -444.8350788756713
Episode * 43 * Avg Reward is ==> -409.85293165991334
Episode * 44 * Avg Reward is ==> -390.83984710631546
Episode * 45 * Avg Reward is ==> -360.88156865913675
Episode * 46 * Avg Reward is ==> -325.26685315168595
Episode * 47 * Avg Reward is ==> -290.2315644399411
Episode * 48 * Avg Reward is ==> -268.0351126010609
Episode * 49 * Avg Reward is ==> -247.8952699063706
Episode * 50 * Avg Reward is ==> -222.99123461788048
Episode * 51 * Avg Reward is ==> -209.0830401020491
Episode * 52 * Avg Reward is ==> -205.65143423678765
Episode * 53 * Avg Reward is ==> -201.8910585767988
Episode * 54 * Avg Reward is ==> -192.18560466037357
Episode * 55 * Avg Reward is ==> -189.43475813660137
Episode * 56 * Avg Reward is ==> -191.92700535454787
Episode * 57 * Avg Reward is ==> -188.5196218645745
Episode * 58 * Avg Reward is ==> -188.17872234729674
Episode * 59 * Avg Reward is ==> -167.33043921566485
Episode * 60 * Avg Reward is ==> -165.01361185173954
Episode * 61 * Avg Reward is ==> -164.5316658073024
Episode * 62 * Avg Reward is ==> -164.4025677076815
Episode * 63 * Avg Reward is ==> -167.27842005634784
Episode * 64 * Avg Reward is ==> -167.12049955654845
Episode * 65 * Avg Reward is ==> -170.02761731078783
Episode * 66 * Avg Reward is ==> -167.56039601863873
Episode * 67 * Avg Reward is ==> -164.60482495249738
Episode * 68 * Avg Reward is ==> -167.45278232469394
Episode * 69 * Avg Reward is ==> -167.42407364484592
Episode * 70 * Avg Reward is ==> -167.57794933965346
Episode * 71 * Avg Reward is ==> -170.6408611483338
Episode * 72 * Avg Reward is ==> -163.96954092530822
Episode * 73 * Avg Reward is ==> -160.82007525469245
Episode * 74 * Avg Reward is ==> -158.38239222565778
Episode * 75 * Avg Reward is ==> -158.3554729720654
Episode * 76 * Avg Reward is ==> -158.51036948298994
Episode * 77 * Avg Reward is ==> -158.68906473090686
Episode * 78 * Avg Reward is ==> -164.60260866654318
Episode * 79 * Avg Reward is ==> -161.5493472156026
```



```

Episode * 80 * Avg Reward is ==> -152.48077012719403
Episode * 81 * Avg Reward is ==> -149.52532010375975
Episode * 82 * Avg Reward is ==> -149.61942419730423
Episode * 83 * Avg Reward is ==> -149.82443455067468
Episode * 84 * Avg Reward is ==> -149.80009937226978
Episode * 85 * Avg Reward is ==> -144.51659331262107
Episode * 86 * Avg Reward is ==> -150.7545561142967
Episode * 87 * Avg Reward is ==> -153.84772667131307
Episode * 88 * Avg Reward is ==> -151.35200443047225
Episode * 89 * Avg Reward is ==> -148.30392250041828
Episode * 90 * Avg Reward is ==> -151.33886235855053
Episode * 91 * Avg Reward is ==> -151.153096135589
Episode * 92 * Avg Reward is ==> -151.19626034791332
Episode * 93 * Avg Reward is ==> -151.15870791946685
Episode * 94 * Avg Reward is ==> -154.2673372216281
Episode * 95 * Avg Reward is ==> -150.40737651480134
Episode * 96 * Avg Reward is ==> -147.7969116731913
Episode * 97 * Avg Reward is ==> -147.88640802454557
Episode * 98 * Avg Reward is ==> -144.88997165191319
Episode * 99 * Avg Reward is ==> -142.22158276699662

```



If training proceeds correctly, the average episodic reward will increase with time.

Feel free to try different learning rates, `tau` values, and architectures for the Actor and Critic networks.

The Inverted Pendulum problem has low complexity, but DDPG work great on many other problems.

Another great environment to try this on is `LunarLandingContinuous-v2`, but it will take more episodes to obtain good results.

```

# Save the weights
actor_model.save_weights("pendulum_actor.h5")
critic_model.save_weights("pendulum_critic.h5")

target_actor.save_weights("pendulum_target_actor.h5")
target_critic.save_weights("pendulum_target_critic.h5")

```

Before Training:



After 100 episodes:

