

spending the libraries

```
import os
import time
import random
import numpy as np
import matplotlib.pyplot as plt
import pybullet_envs
import gym
import torch
import torch.nn as nn
from torch.nn import functional as F
import gym import wrappers
from torch.autograd import Variable
from collections import deque
```

12 | 18.5
12
6.5

Step 1: We initialize the Experience Replay memory

class ReplayBuffer (object):

Step 1: We initialize the experience replay memory.

```
class ReplayBuffer(object):  
    def __init__(self, max_size = 1e6):
```

```
def init (self, max_size = 100):
```

self.storage = [] Memory

self.max_size = max size

$$\text{Sub. ptr} = 0.$$

Index of different cells: (x_1, s_1, a_1, s_2) (x_2, s_2, a_2, s_3) (x_3, s_3, a_3, s_4) (x_4, s_4, a_4, s_5)

```
def add(self, transition):
```

if len(self.storage) == self.max_size:

self-storage $\text{int}(\text{self}, \text{ptr})$ = transition 2

Self. ptr = (self. ptr + 1) % self. max_size

```

else:
    self.storage.append(transition)

```

def sample(self, batch_size):

```
ind = np.random.randint(0, len(self.storage), size=batch_size)
```

batch_states, batch_next_state, batch_actions, batch_rewards, batch_dones = [], [], [], [], []

for i in ind :

state, next_state, action, reward, done = self.storage[i]

```
batch_states.append(np.array(state, copy=False))
```

```
batch_next_states.append(np.array(next_state, copy=False))
```

```
batch_actions.append(np.array(action, copy=False))
```

batch_rewards.append(np.array(reward, copy=False))

both done: append (np.array (done, copy = False))

return np.array(batch_data), np.array(batch_next_data), np.array(batch_labels)

np.array (batch_size, num_channels, num_filters, num_filters) . reshape (-1, 1), np.array (batch_size, num_filters, num_filters)

horizontal 15 array

Step 2: We build one neural network for the Actor model and one neural network for the Actor target.

Class Actor (nn.Module):

```
def __init__(self, state_dim, action_dim, max_action):
    super(Actor, self).__init__()
    self.layer_1 = nn.Linear(state_dim, 400)
    self.layer_2 = nn.Linear(400, 300)
    self.layer_3 = nn.Linear(300, action_dim)
    self.max_action = max_action

def forward(self, x):
    x = F.relu(self.layer_1(x))
    x = F.relu(self.layer_2(x))
    x = self.layer_3(x) * torch.tanh(self.max_action)
    return x
```

Twist Delayed DDPG (TD3)



We build two neural networks for the two Critic models and two neural networks for the two Critic targets.
Critic (nn.Module):

```
def __init__(self, state_dim, action_dim):
    super(Critic, self).__init__()
    # Defining the first Critic Neural Network
    self.layer_1 = nn.Linear(state_dim + action_dim, 400)
    self.layer_2 = nn.Linear(400, 300)
    self.layer_3 = nn.Linear(300, 1)  # Output layer
    # Defining the second Critic neural network.
    self.layer_4 = nn.Linear(state_dim + action_dim, 400)
    self.layer_5 = nn.Linear(400, 300)
    self.layer_6 = nn.Linear(300, 1)  # Output layer

    def forward(self, x, u):
        # Propagate signal from input layer to output layers
        xu = torch.cat([x, u], 1)
        # Forward-Propagation on the first Critic Neural Network
        x1 = F.relu(self.layer_1(xu))
        x1 = F.relu(self.layer_2(x1))
        x1 = self.layer_3(x1)
        # Forward-Propagation on the second Critic neural network.
        x2 = F.relu(self.layer_4(xu))
        x2 = F.relu(self.layer_5(x2))
        x2 = self.layer_6(x2)
        return x1, x2

def Q1(self, x, u):
    xu = torch.cat([x, u], 1)
    x1 = F.relu(self.layer_1(xu))
    x1 = F.relu(self.layer_2(x1))
    x1 = self.layer_3(x1)
    return x1
```

Annotations:

- nn.Module* (written vertically on the left)
- one-to-one* (written vertically on the left, next to the forward method)
- state action* (written above the `xu` variable)
- ReLU* (written above the first `F.relu` call)
- These are going to concatenate action-state* (written next to the `torch.cat` call)
- output loss* (written vertically on the left, next to the second Critic network)
- critic method* (written vertically on the left, next to the `Q1` method)
- if we want to find first output - use the first Critic model* (written on the right, next to the `Q1` method)

selecting the device.
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

Building the whole Training Process into a class
class TD3 Object:

```
def __init__(self, state_dim, action_dim, max_action):  
    self.actor = Actor(state_dim, action_dim, max_action).to(device)  
    self.actor_target = Actor(state_dim, action_dim, max_action).to(device)  
    self.actor_target.load_state_dict(self.actor.state_dict())  
    self.actor_optimizer = torch.optim.Adam(self.actor.parameters())  
    self.critic = Critic(state_dim, action_dim).to(device)  
    self.critic_target = Critic(state_dim, action_dim).to(device)  
    self.critic_target.load_state_dict(self.critic.state_dict())  
    self.critic_optimizer = torch.optim.Adam(self.critic.parameters())  
    self.max_action = max_action
```

```
def select_action(self, state):  
    state = torch.Tensor(state.reshape(1, -1)).to(device)  
    return self.actor(state).cpu().data.numpy().flatten() # convert to list  
# no. of transitions in each batch is back to numpy for clipping and logging
```

```
def train(self, replay_buffer, iterations, batch_size=100, discount=0.99, tau=0.005,  
        policy_noise=0.2, noise_clip=0.5, policy_freq=2):
```

Exploration (standard deviation in our parameters)
+ update of parameters in policy
+ Optimised to train
Frequency of the delay
Formula of Target (used) & default

for it in range(iterations):

```
# Step 4: We sample a batch of transitions (s, s', a, r) from the memory  
batch_states, batch_next_states, batch_actions, batch_rewards, batch_dones = replay_buffer.sample(batch_size)  
state = torch.Tensor(batch_states).to(device)  
next_state = torch.Tensor(batch_next_states).to(device)  
action = torch.Tensor(batch_actions).to(device)  
reward = torch.Tensor(batch_rewards).to(device)  
done = torch.Tensor(batch_dones).to(device)
```

```
# Step 5: From the next state s', the Actor target plays the next action a'  
next_action = self.actor_target(next_state)
```


Step 6: We add Gaussian noise to the next action a' and we clamp it in a range of values supported by the environment.
 $noise = torch.randn(1, policy_noise_std)$
 $noise = noise.clamp(-noise_clip, noise_clip)$
 $next_action = (next_action + noise).clamp(self.min_action, self.max_action)$
clip the noise

Step 7: The two Critic targets take each the couple (s', a') as input and return two Q values $Q_{t1}(s', a')$ and $Q_{t2}(s', a')$ as outputs.

$target_Q1, target_Q2 = self.critic_target(next_state, next_action)$

Step 8: We keep the minimum of these two Q-values: $\min(Q_{t1}, Q_{t2})$
 $target_Q = torch.min(target_Q1, target_Q2)$

Step 9: We get the final target of the two Critic models, which is:
 $Q_t = r + \gamma * \min(Q_{t1}, Q_{t2})$, where γ is the discount factor.
 $target_Q = reward + ((1 - done) * discount * target_Q).detach()$
interpolate done to couple *minimum for value* *detach to avoid graph*

Step 10: The two Critic models take each the sample (s, a) as input and return two Q-values $Q_1(s, a)$ and $Q_2(s, a)$ as outputs.
 $current_Q1, current_Q2 = self.critic(state, action)$

Step 11: We compute the loss coming from the two Critic models:
 $Critic_Loss = MSELoss(Q1(s, a), Q_t) + MSELoss(Q2(s, a), Q_t)$
 $critic_loss = F.mse_loss(current_Q1, target_Q) + F.mse_loss(current_Q2, target_Q)$

Step 12: We backpropagate the Critic loss and update the parameters of the two Critic models with a SGD optimizer.
 $self.critic_optimizer.zero_grad()$ *Initialize gradient to 0s*
 $critic_loss.backward()$ *Backpropagate*
 $self.critic_optimizer.step()$ *Update weight of both Critic models*

Step 13: Once every two iterations, we update our Actor model by performing gradient ascent on the output of the first Critic model.

It is to policy freq = 0.5
 $actor_loss = self.critic_Q1(state, self.actor(state)).mean()$
 $self.actor_optimizer.zero_grad()$ *Backpropagate*
 $self.actor_optimizer.step()$ *Update weight of Actor model*

we perform gradient ascent to increase Q values

p 14: still once every two iterations, we update the weights of the Actor target by Polyak averaging

For param, target_param in zip(self.actor.parameters(), self.actor_target.parameters()):
target_param.data.copy_ = tau * param.data + (1 - tau) * target_param.data

ep 15: still once every two iterations, we update the weights of the Critic target by Polyak averaging

For param, target_param in zip(self.critic.parameters(), self.critic_target.parameters()):
target_param.data.copy_ = tau * param.data + (1 - tau) * target_param.data

EMOTION AI

(CONFUSION MATRIX)