

# 2020 年臺灣國際科學展覽會

## 研究報告

區別：

科別：電腦科學與資訊工程科

作品名稱：機器人自我調整強化學習演算法

關鍵詞：強化學習、自回饋、機器人

編號：

## 目錄

目錄.....	1
---------	---

壹、	研究動機 .....	4
貳、	研究目的 .....	5
參、	文獻探討 .....	5
肆、	研究方法 .....	6
1.	理論背景.....	6
2.	架構構想.....	7
3.	DDPG 簡述 .....	8
4.	算法設計.....	9
5.	實驗環境.....	11
伍、	研究過程 .....	14
1.	正常機體訓練.....	14
陸、	研究規劃 .....	17
柒、	參考資料 .....	18
1.	英文文獻.....	18
2.	網路資源.....	18
捌、	附錄 .....	19
1.	Python 程式碼.....	19

# 摘要

深度強化學習的發展日趨成熟，已經能運用在機器人控制上。隨著機器人的功能多樣化，結構、零件的複雜性大幅提升，機體受損時，若需要人工為各種可能的情形設計一套應對方案，會非常耗時且不切實際。因此若能透過強化學習與環境進行即時互動的特性，使機器人在意外發生時自行學習適應，將能應用在眾多場景中如：遠端遙控、無監督自動運作等。

本研究探討當使用強化學習控制機器人時，如何即時適應機體遭受的損傷。透過模仿人類受傷時可透過去經驗快速適應新情境的機制，提出一套完整的訓練演算法，利用預先生成受損經驗、評估經驗價值等機制，使機器人在實際運作且受損時可以免於重新訓練模型之災，在短時間內適應新情況。

預期能透過模擬四足機器人，逐步測試、分析模型成效，最後使機器人不需依靠額外幫助便在短時間內適應損傷。

# 壹、研究動機

組裝、操控機器人一直是我平日的一大樂趣，二足機器人格鬥賽更是令我熱血沸騰。然而對於各種不同的機體，傳統上皆需人為編寫動作，儘管只是稍微修改結構也可能花上數小時甚至數天來修改動作，因此便有透過強化學習來讓機器人自行學會做出各種動作。



圖 1：DeepMind 訓練 AI 學會跑步(2017)

然而當機器人的零件故障或結構稍有改變時，訓練好的模型極有可能需要重新花費時間訓練。若與人類受傷時做比較，除了人類具有強大的推理能力外，若先前曾有類似的受傷經驗，便能更快速的適應。因此想到了若在訓練過程中讓機器人擁有損傷經驗，或許就能用來在未來的相似情況中達到幫助。

若能擁有適應損傷的能力，機器人便能在風險較大的環境中執行任務，例如遠距操作時，由於無法及時維修，機器便需要在有缺陷的情況下完成任務。

## 貳、研究目的

以往的算法中模型是固定，即模型針對一個機體訓練好後不容易處理不同情況，因此機體受損時很難即時適應，儘管有研究透過搜尋法，預先建立動作－價值對應關係，當機體複雜性增加時，仍會出現效率低落的現象。

因此本次研究希望運用人類在適應損傷時能活用經驗的能力來建立一套算法，使其能在機體損傷的狀況下利用過去經驗修正，避免重新訓練並以較短的時間適應。

## 參、文獻探討

Cully, Clune, Tarapore, & Mouret (2015) 提出利用 Performance Map 建立特徵與模型表現對映表的算法 (以下稱之為 Map 算法)。當機體受損時，利用 Bayesian Optimization 在動作空間中搜索出新策略。例如使用六足蛛型機器人時，使用腿部與地面接觸時間建立 Performance Map、多節機器手臂則使用關節角度等。

此方法需要對不同機型設計不同建立 Performance Map 的特徵，無法普遍適用在不同機型上。且關節要以何種方式做出特徵代表的狀態需要事先設計。

本研究提出使用強化學習(Reinforcement Learning，簡稱 RL)，預先生成損傷經驗，待運作時受到損傷後，提取損傷經驗，並結合實際情況訓練的算法 (以下稱之為經驗適應算法)。經驗適應算法與 Map 算法比較有以下優勢：

表 1：經驗適應算法與 Map 算法之比較

	經驗適應算法	Map 算法
適用性	可直接套用於任何機型	需為不同機型選擇特徵
人工設計動作	所有關節動作皆以 RL 直接控制，無須人工介入	若特徵的形成由多個關節參與，則需人工編寫動作

設計經驗適應算法主要有兩個問題需要解決：

1. 如何生成經驗並依據損傷經驗的價值保存樣本。
2. 適應時如何結合現實狀況與經驗有效訓練模型。

# 肆、研究方法

## 1. 理論背景

強化學習(Reinforcement Learning, 簡稱 RL)是機器學習中的一個領域，與一般機器學習不同的是，其著重在如何與環境的互動以取得最大利益。由於 RL 具有普適性，因此在各種領域如博弈論、控制論、群體智能等皆有研究。在控制系統中，強化學習成功在複雜、多變的情境如機器人控制中取得突出的成果。

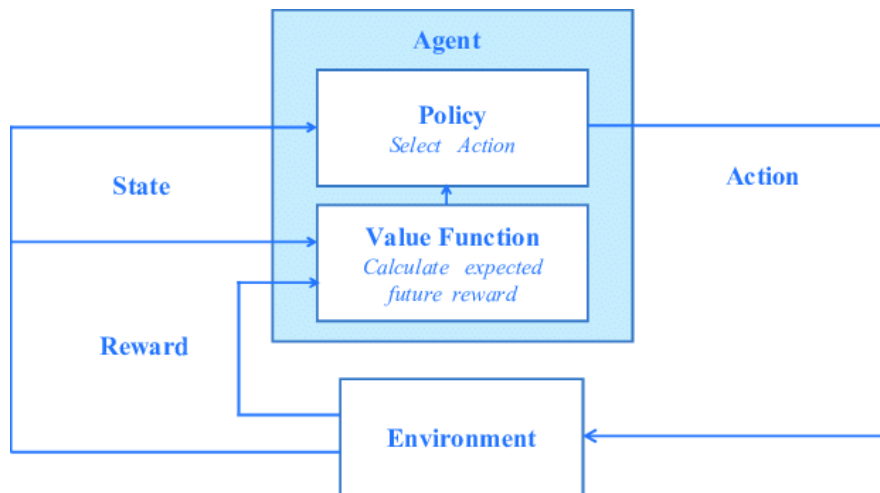


圖 2：簡易強化學習算法架構

一個簡易強化學習模型包含：代理(Agent)、環境(Environment)。在每輪的學習中，Agent 會執行一個動作(Action)，然後從 Environment 得到獎賞(Reward)與狀態(State)，接著透過價值函數(Value Function)估計未來回報，最後經策略(Policy)執行下一個 Action。

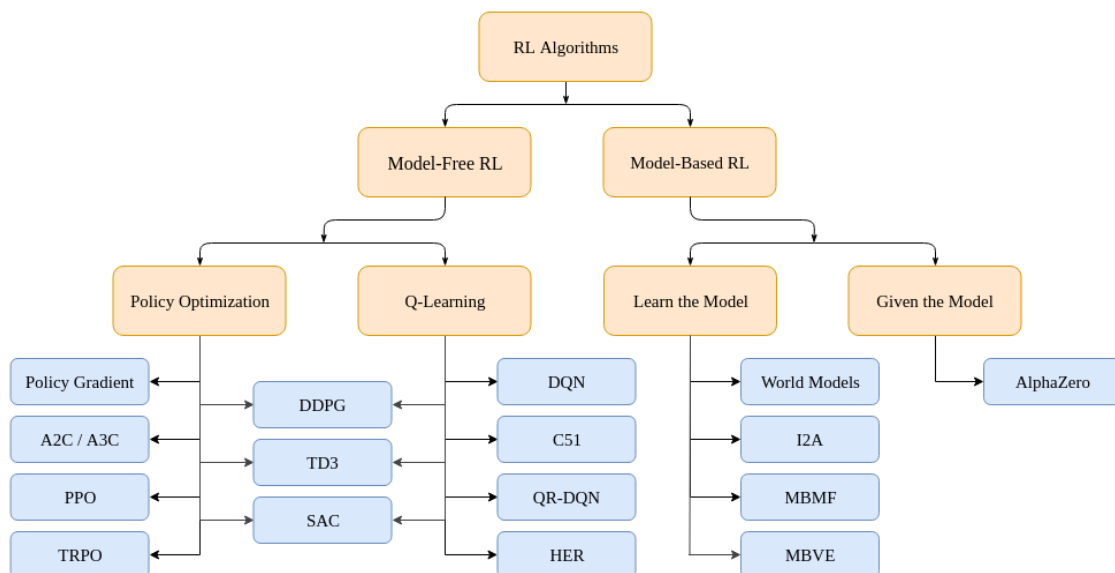


圖 3：強化學習算法生態圖

RL 算法主要可分為兩類：無關模型(Model-Free)、基於模型(Model-Based)

表 2：Model-free/Model-based 比較表

	Model-free	Model-based
說明	算法本身不了解環境狀態可能如何轉移，透過長期與環境互動的收益，仍可以評估與預測不同動作的回報	在學習的過程可觀察與學習環境在不同動作下的變化機制，預測環境狀態如何轉移
優點	由於能在未知環境中學習，因此具有更好的普適性	充分了解環境，可利用動態規劃算法求出較好的動作
缺點	學習僅能透過探索評估動作好壞，通常需較多時間收斂	大多數情況下無法提供算法對於環境的精確描述導致算法普適性較低

實際上面對複雜多變的環境難以用規劃法求解，因此多數 RL 算法 Model-free 的。而 Model-free 算法又可依據優化的對象的再區分為兩類：基於值(Value-based)、基於策略(Policy-based)

表 3：Value-based/Policy-based 比較表

	Value-based	Policy-based
說明	學習如何正確評估動作的好壞，也就是優化 Value Function	學習如何做出最好的動作，也就是優化 Policy
舉例	DQN、SARSA、QR-DQN	Policy Gradient、PPO、TRPO

## 2. 架構構想

### 2.1. 算法基礎

本研究採用 DDPG(Deep Deterministic Policy Gradient)算法作為基礎，原因在於 DDPG 混和 Value-based 與 Policy-based，能處理評價與選擇動作的複雜性；且 DDPG 能處理連續動作空間，適合馬達的運作。

### 2.2. 經驗生成

在 RL 中，為了增加訓練效率與模型泛化性，使用經驗回放(Experience Replay)機制，將模型經歷過的經驗儲存起來作為訓練樣本。因此本研究利用此機制，在訓練階段預先在受損狀態下進行定量訓練，獲取相關樣本，以利在實際受損時幫助訓練。

### 2.3. 適應訓練

通常實際運作時的環境不會與訓練完全相同，受損情況也會有差異，因此篩選出損傷經驗後，機器仍需同時進行探索，補充實際狀況的樣本。

### 3. DDPG 簡述

#### 3.1. Actor-Critic

DDPG 採用 Actor-Critic 來達到同時學習 Policy 與 Value Function：

##### 3.1.1. Critic

神經網路，作為 Value Function 使用，輸入為 Action  $a$  與 State  $s$ ，輸出為  $Q(s, a; \omega)$ ， $\omega$  為網路參數。Critic 用來評估未來回報：

$$Q(s_t, a_t; \omega) = E \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} \middle| s_t, a_t \right] = E[r_t + \gamma Q(s_{t+1}, \mu(s_{t+1}; \theta)) | s_t, a_t] \quad (1)$$

$$Q^*(s_t, a_t) = \max_{\omega} Q(s_t, a_t; \omega) = E[r_t + \gamma Q^*(s_{t+1}, \mu(s_{t+1}; \theta)) | s_t, a_t]$$

其中  $0 < \gamma < 1$  為折扣值，表示對於越遙遠的未來，重要性越低。

訓練時，根據  $Q(s, a; \omega)$  的定義來設定更新目標並計算損失：

$$\begin{aligned} target_t &= r_t + \gamma Q(s_{t+1}, \mu(s_{t+1}; \theta^-); \omega^-) \\ loss &= (target_t - Q(s_t, a_t; \omega))^2 \end{aligned} \quad (2)$$

##### 3.1.2. Actor

神經網路，作為 Policy 使用，輸入為狀態  $s$ ，輸出為  $\mu(s; \theta)$ ， $\theta$  為網路參數。

透過策略梯度更新 Actor，令  $J(\theta)$  為評價策略的函數：

$$J(\theta) = \mathbb{E}_s [Q(s, \mu(s; \theta))] \quad (3)$$

策略梯度即為：

$$\Delta_{\theta} J(\theta) = \mathbb{E}_s [\Delta_a Q(s, a) |_{a=\mu(s)} \Delta_{\theta} \mu(s; \theta)] \quad (4)$$

兩個網路的交互關係為： $t$  時刻時，Actor 根據  $s_t$  作出決策  $a_t = \mu(s_t; \theta)$ ，接收 Reward  $r_{t+1}$  後，用  $r_{t+1}$  更新 Critic，然後以 Critic 給出的  $Q(s_t, a_t; \omega)$  更新 Actor。

由於更新網路時，若使用網路本身來計算目標，會導致以自己為目標更新。因此 Actor 與 Critic 會另外使用 Target Actor 與 Target Critic 作為更新對象，網路參數分別為  $\theta^-$ ,  $\omega^-$ 。

Target 初始與原網路相同，訓練採 Soft-Update：

$$\begin{cases} \theta^- \leftarrow \tau \theta + (1 - \tau) \theta^- \\ \omega^- \leftarrow \tau \omega + (1 - \tau) \omega^- \end{cases} \text{ with } \tau \ll 1 \quad (5)$$

##### 3.1.3. Experience Replay

令經驗庫為  $\mathcal{D}$ ，其中儲存樣本  $(s_i, a_i, r_i, s_{i+1})$ ，容量限制為  $\mathcal{C}$ 。每個時間步從  $\mathcal{D}$  取出一個 batch 的樣本  $(s_j, a_j, r_j, s_{j+1})$  更新 Actor 與 Critic。

$$L(\omega) = \frac{1}{N} \sum_{k=1}^N (target_k - Q(s_t, a_t; \omega))^2 \quad (6)$$

對於 Actor，沿策略梯度方向更新：

$$\Delta_{\theta} J(\theta) = \frac{1}{N} \sum_{k=1}^N \Delta_a Q(s, a) |_{s=s_k, a=\mu(s_k)} \Delta_{\theta} \mu(s; \theta) |_{s=s_k} \quad (7)$$



然後將當前樣本  $(s_t, a_t, r_t, s_{t+1})$  存入  $\mathcal{D}$ ，並使用時間差分誤差(TD-error)評估樣本值得不值得學習與保留：

## 4. 算法設計

#### 4.1. 訓練階段

圖 4：訓練階段流程

圖 5：Train 模塊流程

#### 4.1.3. 評估適應難度(Evaluate value of experience)

評估損傷的適應難度，低則保留少量樣本；高則保留較多樣本。若訓練後平均回報高、回報波動幅度低則適應難度低；若平均回報低、回報波動幅度高則適應難度低。

波動幅度計算時，先透過計算回報的時間差分序列消除原序列趨勢，再計算整體標準差  $\sigma_j$ ，令 *episode* 總數為  $EP$ ，每個 *episode* 總回報  $R_j^{ep}$ ：

$$\sigma_j = \sqrt{\frac{1}{EP-1} \sum_{ep=2}^{EP} (\Delta R_j^{ep} - \mu(\Delta R_j))^2} \quad (9)$$

$$\Delta R_j^{ep} = R_j^{ep} - R_j^{ep-1}$$

$$\mu(\Delta R_j) = \frac{1}{EP-1} (R_j^{EP} - R_j^1)$$

令平均回報為  $\bar{R}_j$ ，回報可為負值，故減去最小值後再計算其指標  $q_j$ ，令 *joint* 總數為  $JN$ ：

$$q_j = e^{-\frac{\bar{R}_j - \min \bar{R}}{\mu(\bar{R}) - \min \bar{R}}} \quad (10)$$

$$\mu(\bar{R}) = \sum_{j=1}^{JN} \bar{R}_j$$

適應難度  $v_j = q_j \sigma_j$

#### 4.1.4. 分配樣本

依據適應難度分配樣本，令  $c_j$  為鎖定 *joint<sub>j</sub>* 的損傷樣本數：

$$c_j = \frac{v_j}{\sum_{i=0}^{JN} v_i} \mathcal{C} \quad (11)$$

### 4.2. 運作階段

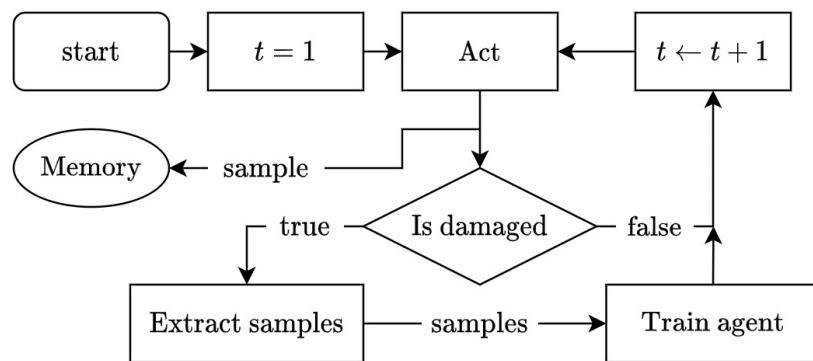


圖 6：運作流程

#### 4.2.1. 執行(Act)

執行任務並儲存運作時的樣本。

#### 4.2.2. 提取樣本(Extract samples)

若機體受損，從經驗庫提取損傷用於訓練 Agent 適應。

## 5. 實驗環境

利用 CoppeliaSim 的遠端 API 給 Python 語言進行操作，再用 TensorFlow 建立模型與訓練。

表 4：實驗環境

作業系統	Windows 10 Home
雲端運算環境	TWCC 臺灣 AI 雲
程式語言	Python 3.6.8
模型框架	TensorFlow 2.0
機器人模擬器	CoppeliaSim V4.0.0

### 5.1. 機體設計

考慮本次研究希望觀察正常狀態至損傷狀態的適應情形，因此挑選容易訓練、損傷影響大的四足機器人。共有 8 個關節：

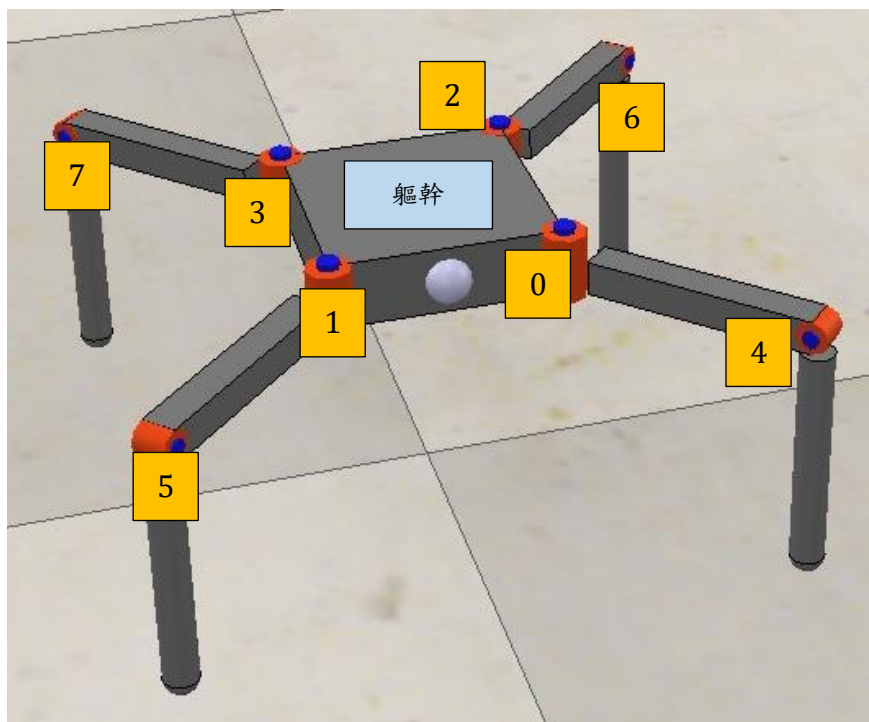


圖 7：模擬用四足機器人

表 5：機體狀態、動作資訊

	維度	資料
State	25	軀幹：座標 $x, y, z$ 、速度 $v_x, v_y, v_z$ 、傾角 $\psi, \theta, \phi$ 馬達：角度 $q_1, q_2, \dots, q_8$ 、角速度 $\omega_1, \omega_2, \dots, \omega_8$
Action	8	馬達：輸出力矩 $\tau_1, \tau_2, \dots, \tau_8$

## 5.2. 場景設計

讓機器人沿著  $x$  正方向移動至目標處，如圖：

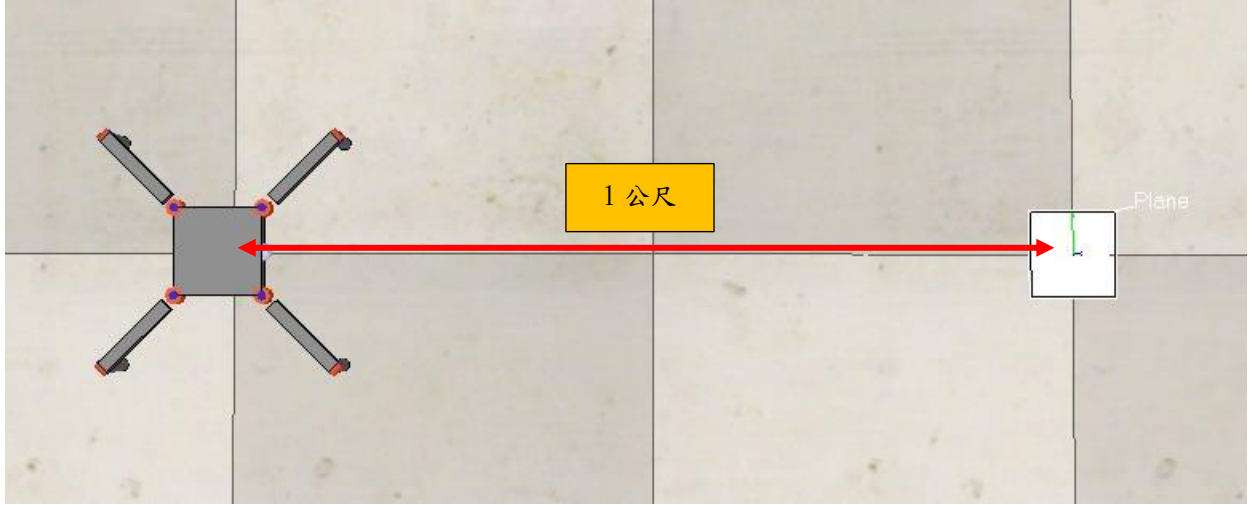


圖 8：Coppeliasim 模擬訓練場景

## 5.3. Reward Function

實驗中的 Reward Function 期望使機器人達成以下目標：

### 5.3.1. 行走平穩

向前移動，故  $v_x$  應保持正值；盡量沿著直線移動，故  $y$  應接近零；減少  $y$  方向上的晃動，故  $v_z$  應接近零。得出機體姿態對 Reward 的貢獻值：

$$w_x v_x - (w_y y^2 + w_z v_z^2) \quad (12)$$

其中  $w_x, w_y, w_z$  為正權重。

### 5.3.2. 動作精簡

希望使機器人能以較有效率的輸出達成目標，因此對馬達總輸出加入懲罰項：

$$-w_\tau \sum_{j=1}^{JN} |\tau_j| \quad (13)$$

其中  $w_\tau$  為正權重。

$$r_t = w_x v_x - (w_y y^2 + w_z v_z^2) - w_\tau \sum_{j=1}^{JN} |\tau_j| \quad (14)$$

#### 5.4. Agent 神經網路結構

Actor 輸入 State，經兩層隱藏輸出：

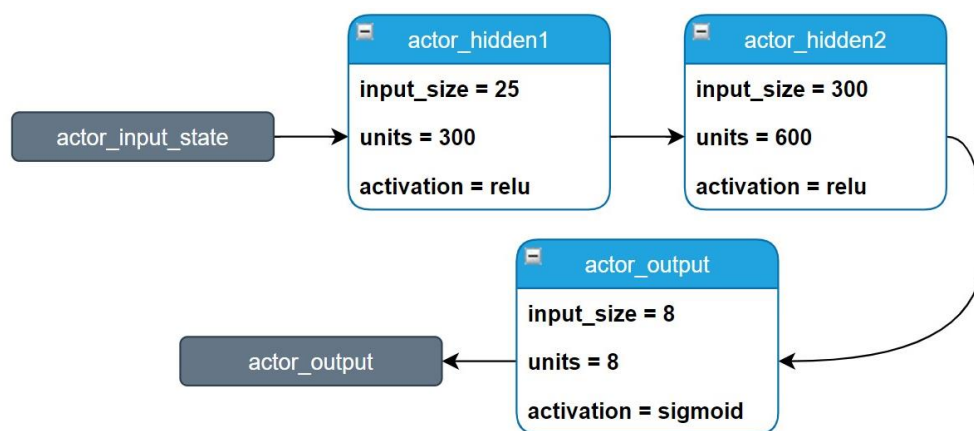


圖 9：Actor 神經網路結構

Critic 輸入 State 與 Action，分別經一層隱藏層後加總，再經一層隱藏輸出：

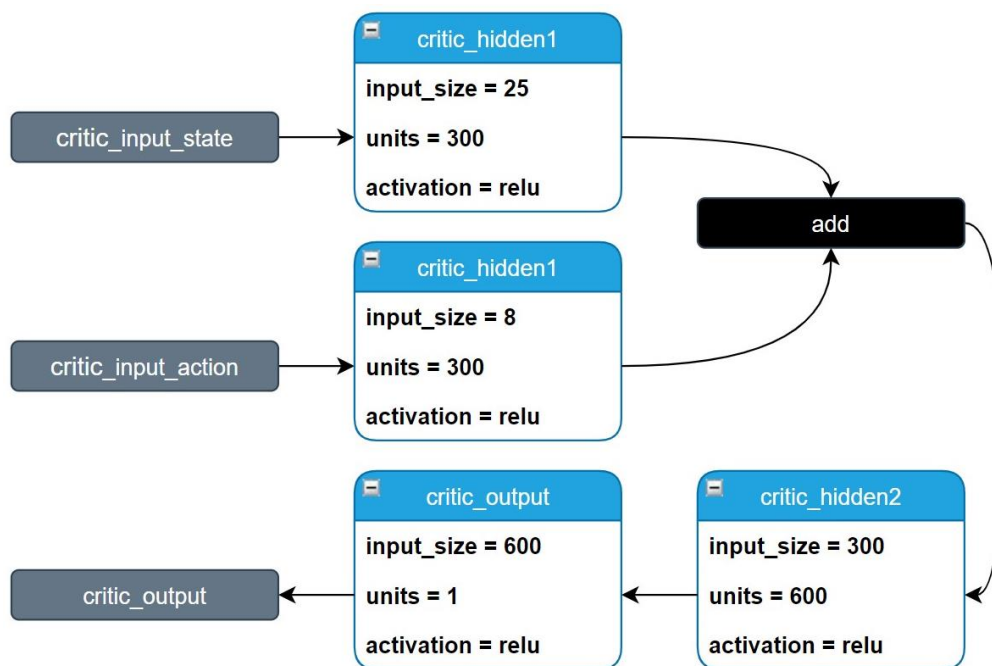


圖 10：Critic 神經網路結構

## 伍、研究過程

### 1. 正常機體訓練

影響 RL 模型表現的最大要素在於 Reward Function 的設計，以下為對 Reward Function 進行調整的過程：

#### 1.1. Reward Function 測試

表 6：基礎參數設置

模型參數	Network Size		$\gamma$	$\tau$	Memory Size
	Actor	Critic	0.99	0.001	1000000
	(300, 600)	(300, 600)			
訓練參數	Learning Rate		Batch Size	Episodes	Steps/Episodes
	Actor	Critic	128	200	200
	0.001	0.01			

##### 1.1.1. 測試 1.1.1

表 7：測試 1.1.1 參數設置

$w_x$	$w_y$	$w_z$	$w_t$
20	15	5	0.05

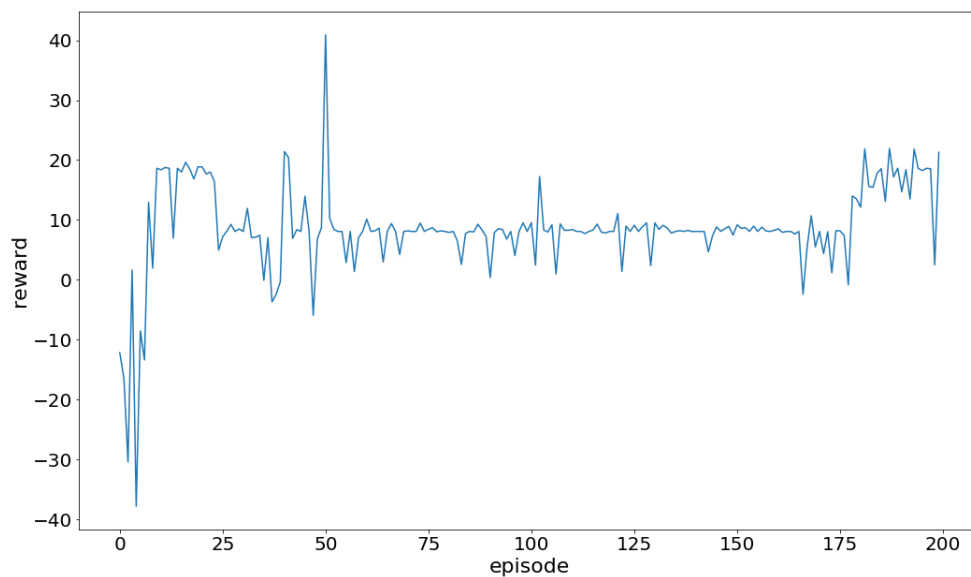


圖 11：測試 1.1.1 結果

由此次結果注意到模型發現保持不動能維持 Reward 在微小正值，因此放棄行動，儘管在 Action 上加入了雜訊來增加探索，效果仍然有限。

### 1.1.2. 測試 1.1.2

根據測試 1.1.1，發現應該將 Reward 修改為：

$$r_t = w_x(v_x - v_{min}) - (w_y y^2 + w_z v_z^2) - w_\tau \sum_{j=1}^{JN} |\tau_j| \quad (15)$$

透過設置最低速度，確保模型能學習加速。

Reward Function 參數設置延續測試 1.1.1。

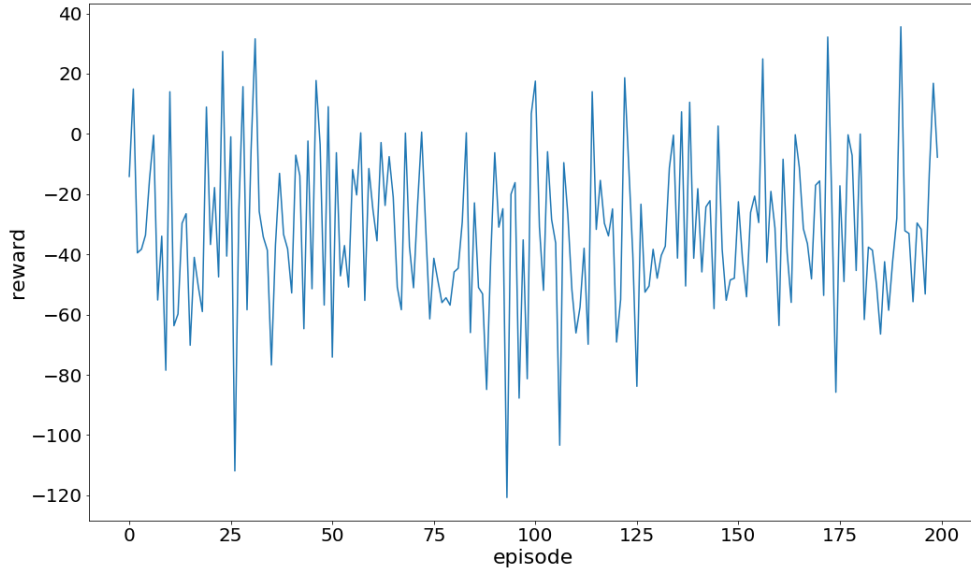


圖 12：測試 1.1.2 結果

不再陷入到靜止不動的狀態。

## 1.2. Learning Rate 調整

### 1.2.1. 測試 1.2.1

模型中兩個神經網路有不同的學習率，觀察測試 1.1.2 結果，發現可能是 Learning Rate 需要調整，因此延續測試 1.1.2，紀錄 Actor 與 Critic 每個 Episode 中的神經元平均梯度大小：

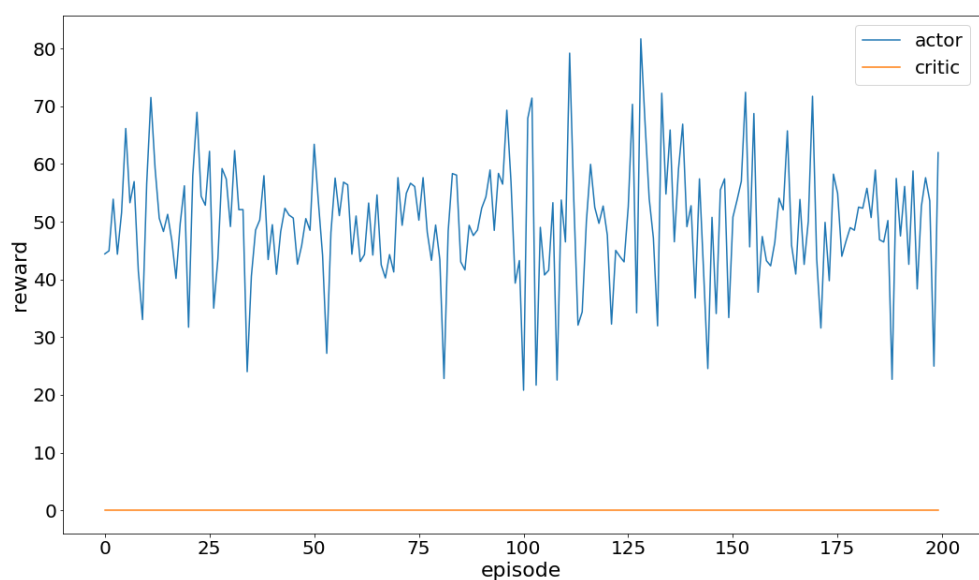


圖 13：測試 1.2.1 神經元平均梯度量值

Actor 平均梯度量值乘上學習率後數量級約為  $10^{-2}$ ，在可接受範圍。然而 Critic 平均梯度量值乘上學習率後數量級約為  $10^{-6}$ ，收斂速度極慢。因此設置新的學習率為：

表 8：修正後學習率

Actor	Critic
0.001	5

### 1.2.2. 測試 1.2.2

延續測試 1.2.1，利用修正過的學習率重新訓練得：

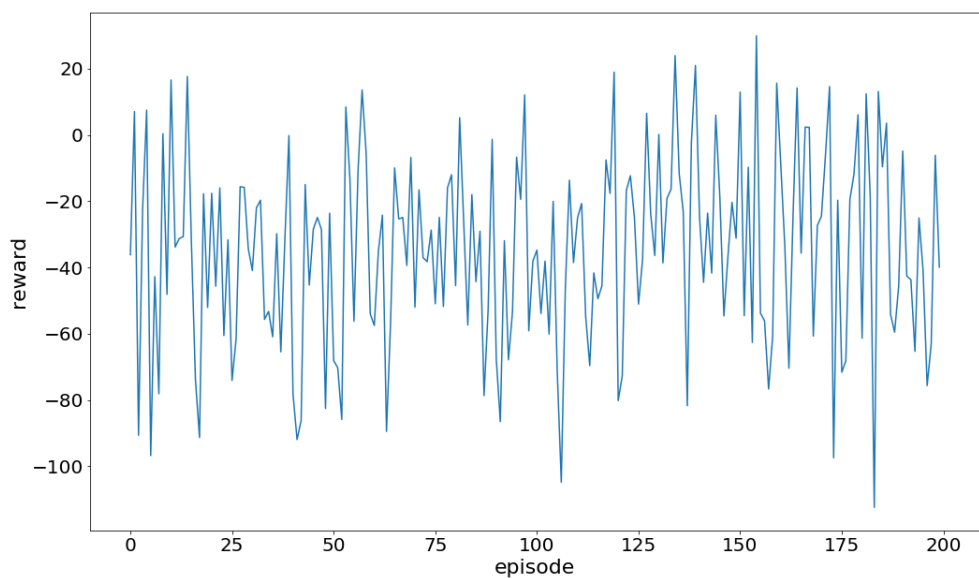


圖 14：測試 1.2.2



## 陸、研究規劃

初步研究結果如下：

表 9：初步研究結果

演算法建構	確立生成經驗之方法
	完成演算法流程設計
	完成模型架構
	建構評估經驗價值之方法
模型訓練	分析 Reward Function 之影響並修正
	分析 Learning Rate 之影響並修正

研究階段規劃如下：

表 10：研究階段規劃

	階段研究內容
正常訓練	分析神經網路結構之影響並修正
經驗生成	分析生成經驗時各參數之影響並修正
	分析經驗價值評估法之影響並修正
適應訓練	測試適應效果

# 柒、參考資料

## 1. 英文文獻

- [1] Tom Schaul, John Quan, Ioannis Antonoglou, & David Silver, Prioritized Experience Replay, 2015. URL <https://arxiv.org/abs/1511.05952>
- [2] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, & Daan Wierstra, Continuous control with deep reinforcement learning, 2015. URL <https://arxiv.org/abs/1509.02971>
- [3] Shixiang Gu, Ethan Holly, Timothy Lillicrap, & Sergey Levine<sup>1</sup>, Deep Reinforcement Learning for Robotic Manipulation with Asynchronous Off-Policy Updates, 2016. URL <https://arxiv.org/abs/1610.00633>
- [4] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, & Martin Riedmiller, Playing Atari with Deep Reinforcement Learning, 2013. URL <https://arxiv.org/abs/1312.5602>
- [5] Hao Yi Ong, Kevin Chavez, & Augustus Hong, Distributed Deep Q-Learning, 2015. URL <https://arxiv.org/abs/1508.04186>
- [6] Antoine Cully, Jeff Clune, Danesh Tarapore, & Jean-Baptiste Mouret, Robots that can adapt like animals, 2015. URL <https://arxiv.org/abs/1407.3501>

## 2. 網路資源

- [1] <Wikipedia. Reinforcement learning>, URL [https://en.wikipedia.org/wiki/Reinforcement\\_learning](https://en.wikipedia.org/wiki/Reinforcement_learning)

# 捌、附錄

## 1. Python 程式碼

### 1.1. DDPG

```
1. import tensorflow as tf
2. import numpy as np
3.
4. from .networks.actor import Actor
5. from .networks.critic import Critic
6. from utils.memory import Memory
7.
8. from tensorflow.keras.losses import MSE
9. from tensorflow.keras.optimizers import Adam
10.
11. class DDPG(object):
12.
13.     def __init__(self,
14.                  state_shape, action_shape,
15.                  action_range, state_range,
16.                  noise_std=0.05,
17.                  actor_lr=1e-4, critic_lr=1e-3,
18.                  actor_hidden_units=(300, 600), critic_hidden_units=(300, 600),
19.                  batch_size=128, discount=0.99, memory_size=1e6, tau=0.001):
20.
21.         self.discount = discount
22.         self.batch_size = batch_size
23.         self.tau = tau
24.
25.         self.state_shape = state_shape
26.         self.action_shape = action_shape
27.
28.         self.action_range = action_range
29.         self.state_range = state_range
30.
31.         self.noise_std = noise_std
32.
33.         # create actor
34.         self.actor = Actor(state_shape=state_shape, action_shape=action_shape, hidden_units=actor_hidden_units)
35.         self.target_actor = Actor(state_shape=state_shape, action_shape=action_shape, hidden_units=actor_hidden_units)
36.         # create critic
37.         self.critic = Critic(state_shape=state_shape, action_shape=action_shape, hidden_units=critic_hidden_units)
38.         self.target_critic = Critic(state_shape=state_shape, action_shape=action_shape, hidden_units=critic_hidden_units)
39.
40.         self.actor_optimizer = Adam(learning_rate=actor_lr)
41.         self.critic_optimizer = Adam(learning_rate=critic_lr)
42.
43.         self.memory = Memory(int(memory_size), action_shape, state_shape)
44.
45.     def store_transition(self, state, action, reward, next_state, done):
46.
47.         self.memory.append(state, action, reward, next_state, done)
48.
```

```

49.     def step(self, state, apply_noise=True):
50.
51.         action = self.actor(state)
52.         if apply_noise:
53.             action += tf.random.normal(self.action_shape, stddev=self.noise_std)
54.         action = tf.clip_by_value(action, self.action_range[0], self.action_range[1])
55.
56.         return action
57.
58.     def train(self):
59.
60.         if len(self.memory) >= self.batch_size:
61.             states, actions, rewards, next_states, done = self.memory.sample(self.batch_size)
62.
63.             ct = self.train_critic(states, actions, next_states, rewards, done)
64.             at = self.train_actor(states)
65.             self.update_target_models()
66.
67.             return ct, at
68.
69.     def train_critic(self, states, actions, next_states, rewards, done):
70.         # q target
71.         next_actions = self.target_actor(next_states)
72.         next_q = self.target_critic(next_states, next_actions)
73.         q_targets = rewards + (1 - done) * self.discount * next_q
74.
75.         weights = self.critic.get_trainable_weights()
76.         with tf.GradientTape() as tape:
77.             tape.watch(weights)
78.             # compute MSE loss
79.             q_values = self.critic.model([states, actions])
80.             loss = MSE(q_targets, q_values)
81.             # compute gradients
82.             grads = tape.gradient(loss, weights)
83.
84.             tot = 0
85.             for w in grads:
86.                 tot += tf.math.reduce_mean(tf.abs(w))
87.             tot /= len(grads)
88.
89.             self.critic_optimizer.apply_gradients(zip(grads, weights))
90.             return np.asarray(tot)
91.
92.     def train_actor(self, states):
93.
94.         weights = self.actor.get_trainable_weights()
95.         with tf.GradientTape() as tape:
96.             tape.watch(weights)
97.             # compute policy value
98.             actions = self.actor(states)
99.             q_values = self.critic(states, actions)
100.            loss = -tf.math.reduce_mean(q_values)
101.            # compute policy gradients
102.            grads = tape.gradient(loss, weights)
103.
104.            self.actor_optimizer.apply_gradients(zip(grads, weights))
105.            return np.asarray(tot)
106.

```

```

1.     def initialize(self):
2.
3.         actor = self.actor.get_trainable_weights()
4.         target_actor = self.target_actor.get_trainable_weights()
5.         for weight, target_weight in zip(actor, target_actor):
6.             target_weight.assign(weight)
7.
8.         critic = self.critic.get_trainable_weights()
9.         target_critic = self.target_critic.get_trainable_weights()
10.        for weight, target_weight in zip(critic, target_critic):
11.            target_weight.assign(weight)
12.
13.    def update_target_models(self):
14.
15.        actor = self.actor.get_trainable_weights()
16.        target_actor = self.target_actor.get_trainable_weights()
17.        for weight, target_weight in zip(actor, target_actor):
18.            target_weight.assign((1. - self.tau) * target_weight + self.tau * weight)
19.
20.        critic = self.critic.get_trainable_weights()
21.        target_critic = self.target_critic.get_trainable_weights()
22.        for weight, target_weight in zip(critic, target_critic):
23.            target_weight.assign((1. - self.tau) * target_weight + self.tau * weight)

```

## 1.2. Actor

```

1. from tensorflow.keras.models import Model
2. from tensorflow.keras.layers import Dense, Input
3.
4. class Actor(object):
5.
6.     def __init__(self, state_shape, action_shape, hidden_units=(300, 600)):
7.
8.         # store parameters
9.         self.state_shape = state_shape
10.        self.action_shape = action_shape
11.        self.hidden = hidden_units
12.
13.        # generate model
14.        self.model = self.generate_model()
15.
16.    def __call__(self, state):
17.
18.        return self.model(state)
19.
20.    def get_weights(self):
21.
22.        return self.model.weights
23.
24.    def get_trainable_weights(self):
25.
26.        return self.model.trainable_weights
27.
28.    def set_weights(self, weights):
29.
30.        self.model.set_weights(weights)
31.

```

```

1.     def generate_model(self):
2.
3.         # input state
4.         input_layer = Input(shape=self.state_shape,
5.                             name="actor_input_state")
6.         # hidden layer1
7.         layer = Dense(self.hidden[0], activation='relu',
8.                       kernel_initializer='random_normal',
9.                       bias_initializer='zeros',
10.                      name="actor_hedden1")(input_layer)
11.        # hidden layer2
12.        layer = Dense(self.hidden[1], activation='relu',
13.                      kernel_initializer='random_normal',
14.                      bias_initializer='zeros',
15.                      name="actor_hedden2")(layer)
16.        # output layer
17.        output_layer = Dense(self.action_shape[0], activation='sigmoid',
18.                             kernel_initializer='random_normal',
19.                             bias_initializer='zeros',
20.                             name="actor_output")(layer)
21.        # create model
22.        model = Model(inputs=input_layer, outputs=output_layer)
23.
24.        return model

```

### 1.3. Critic

```

1. from tensorflow.keras.layers import Dense, Input, Add
2. from tensorflow.keras.models import Model
3.
4. class Critic(object):
5.
6.     def __init__(self, state_shape, action_shape, hidden_units=(300, 600)):
7.
8.         # store parameters
9.         self.state_shape = state_shape
10.        self.action_shape = action_shape
11.        self.hidden = hidden_units
12.
13.        # generate model
14.        self.model = self.generate_model()
15.
16.    def __call__(self, state, action):
17.
18.        return self.model([state, action])
19.
20.    def get_trainable_weights(self):
21.
22.        return self.model.trainable_weights
23.
24.    def set_weights(self, weights):
25.
26.        self.model.set_weights(weights)
27.
28.    def get_weights(self):
29.
30.        return self.model.weights
31.

```

```

1.     def generate_model(self):
2.
3.         # input state
4.         s_input_layer = Input(shape=self.state_shape,
5.                                name="critic_input_state")
6.         # input action
7.         a_input_layer = Input(shape=self.action_shape,
8.                                name="critic_input_action")
9.         # state hidden layer
10.        s_layer = Dense(self.hidden[0], activation='relu',
11.                        kernel_initializer='random_normal',
12.                        bias_initializer='zeros',
13.                        name="critic_state_hidden_layer")(s_input_layer)
14.        # action hidden layer
15.        a_layer = Dense(self.hidden[0], activation='linear',
16.                        kernel_initializer='random_normal',
17.                        bias_initializer='zeros',
18.                        name="critic_action_hidden_layer")(a_input_layer)
19.        # add state hidden layer and action hidden layer
20.        hidden = Add()([s_layer, a_layer])
21.        # hidden layer2
22.        hidden = Dense(self.hidden[1], activation='relu',
23.                        kernel_initializer='random_normal',
24.                        bias_initializer='zeros',
25.                        name="critic_hidden_layer2")(hidden)
26.        # output layer
27.        output_layer = Dense(1, activation='linear',
28.                              kernel_initializer='random_normal',
29.                              bias_initializer='zeros',
30.                              name="critic_output")(hidden)
31.        # create model
32.        model = Model(inputs=[s_input_layer, a_input_layer],
33.                       outputs=output_layer)
34.
35.        return model

```

#### 1.4. Environment

```

1. import numpy as np
2. from .simAPI import sim
3.
4. _oneshot = sim.simx_opmode_oneshot
5. _blocking = sim.simx_opmode_blocking
6.
7. PI = 3.14159265358979
8. DEG2RAD = PI / 180
9.

```

```

1. class Environment(object):
2.
3.     def __init__(self, state_shape, action_shape,
4.                  dt=0.010, frames=10,
5.                  min_vel = 0.10):
6.
7.         self.state_shape = state_shape
8.         self.action_shape = action_shape
9.
10.        self.dt = dt
11.        self.frames = frames
12.
13.        self.min_vel = min_vel
14.
15.        self.body_name = "QuadSpider"
16.
17.        self.joint_num = 8
18.        self.max_joint_torque = 10
19.        self.max_joint_velocity = 360 * DEG2RAD
20.
21.        self.w_vx = 20
22.        self.w_y = 15
23.        self.w_z = 5
24.
25.        self.w_tor = 0.05
26.
27.    def open(self):
28.
29.        # clear communications
30.        sim.simxFinish(-1)
31.        # create communication
32.        self.client_ID = sim.simxStart('127.0.0.1', 19997, True, True, 5000, 5)
33.
34.        self.stop_sim()
35.
36.        # read body's handle
37.        _, self.body_handle = sim.simxGetObjectHandle(self.client_ID, self.body_name, _blocki
ng)
38.
39.        # read joints' handles
40.        self.joint_handles = []
41.        for i in range(self.joint_num):
42.            _, res = sim.simxGetObjectHandle(self.client_ID, "joint" + str(i), _blocking)
43.            self.joint_handles.append(res)
44.
45.        sim.simxGetPingTime(self.client_ID)
46.
47.    def close(self):
48.
49.        # stop simulation before close communication
50.        self.stop_sim()
51.        # close communication
52.        sim.simxFinish(self.client_ID)
53.
54.        sim.simxGetPingTime(self.client_ID)
55.
56.    def start_sim(self):
57.
58.        # set synchronous
59.        sim.simxSynchronous(self.client_ID, True)
60.        # start simulation
61.        sim.simxStartSimulation(self.client_ID, _oneshot)
62.        # trigger once
63.        sim.simxSynchronousTrigger(self.client_ID)
64.
65.        sim.simxGetPingTime(self.client_ID)
66.

```



```

1.     def stop_sim(self):
2.
3.         sim.simxStopSimulation(self.client_ID, _oneshot)
4.
5.         sim.simxGetPingTime(self.client_ID)
6.
7.     def initialize(self):
8.         for i in range(self.joint_num):
9.             # set the direction of target velocity
10.            sim.simxSetJointTargetVelocity(self.client_ID, self.joint_handles[i],
11.                                           self.max_joint_velocity, _oneshot)
12.            sim.simxSetJointForce(self.client_ID, self.joint_handles[i],
13.                                 self.max_joint_torque, _oneshot)
14.
15.        _, position = sim.simxGetObjectPosition(self.client_ID, self.body_handle, -
16.        1, _blocking)
17.        self.height = position[2]
18.
19.        self.time = 0
20.
21.    def act(self, act):
22.
23.        # pause communication until all commands are sent
24.        sim.simxPauseCommunication(self.client_ID, True)
25.
26.        act = act[0]
27.        # set target position
28.        for i in range(self.joint_num):
29.            tor = (act[i] - 0.5) * 2 * self.max_joint_torque
30.            # set the direction of target velocity
31.            if tor >= 0:
32.                sim.simxSetJointTargetVelocity(self.client_ID, self.joint_handles[i],
33.                                                self.max_joint_velocity, _oneshot)
34.                sim.simxSetJointForce(self.client_ID, self.joint_handles[i],
35.                                     tor, _oneshot)
36.            else:
37.                sim.simxSetJointTargetVelocity(self.client_ID, self.joint_handles[i],
38.                                                -self.max_joint_velocity, _oneshot)
39.                sim.simxSetJointForce(self.client_ID, self.joint_handles[i],
40.                                     -tor, _oneshot)
41.
42.        # restart communication
43.        sim.simxPauseCommunication(self.client_ID, False)
44.        # trigger once
45.        for i in range(self.frames):
46.            sim.simxSynchronousTrigger(self.client_ID)
47.            sim.simxGetPingTime(self.client_ID)
48.            self.time += self.dt
49.
50.    def get_reward(self, velocity, position, torques):
51.
52.        r_vx = self.w_vx * velocity[0]
53.        r_y = self.w_y * position[1] * position[1]
54.        r_z = self.w_z * (position[2] - self.height) * (position[2] - self.height)
55.
56.        r_tor = self.w_tor * np.square(torques).sum()
57.
58.        return r_vx - r_y - r_z - r_tor

```

```

1.     def is_done(self, velocity, position, orientation):
2.
3.         if position[0] >= 3:
4.             return True
5.
6.         if abs(orientation[0]) >= PI / 2:
7.             return True
8.         if abs(orientation[1]) >= PI / 2:
9.             return True
10.        if abs(orientation[2]) >= PI / 2:
11.            return True
12.
13.        if self.time >= 2:
14.            if position[0] < 0:
15.                return True
16.            if abs(position[1]) > 0.2:
17.                return True
18.            if velocity[0] < 0.05:
19.                return True
20.
21.        return False
22.
23.    def get_state(self):
24.
25.        # read body position
26.        _, position = sim.simxGetObjectPosition(self.client_ID, self.body_handle, -
1, _blocking)
27.        position = np.array(position)
28.        # read body velocity
29.        _, velocity, _ = sim.simxGetObjectVelocity(self.client_ID, self.body_handle, _blockin
g)
30.        velocity = np.array(velocity)
31.        # read body orientation
32.        _, orientation = sim.simxGetObjectOrientation(self.client_ID, self.body_handle, -
1, _blocking)
33.        orientation = np.array(orientation)
34.
35.        joint_torques = []
36.        joint_velocities = []
37.        for i in range(self.joint_num):
38.            # read joint force
39.            _, tor = sim.simxGetJointForce(self.client_ID, self.joint_handles[i], _blocking)
40.
41.            # read joint velocity
42.            _, vel = sim.simxGetObjectFloatParameter(self.client_ID, self.joint_handles[i], s
im.sim_jointfloatparam_velocity, _blocking)
43.            vel *= DEG2RAD
44.            # record joint state
45.            joint_torques.append(tor)
46.            joint_velocities.append(vel)
47.
48.        joint_torques = np.array(joint_torques)
49.        joint_velocities = np.array(joint_velocities)
50.
51.        sim.simxGetPingTime(self.client_ID)
52.
53.        return np.concatenate([position, velocity, orientation,
                                joint_torques, joint_velocities]).reshape((1, *self.state_shap
e)), \
54.            self.get_reward(velocity, position, joint_torques), \
55.            self.is_done(velocity, position, orientation)

```