

EXPLORACIÓN DEL DDPG TD3: UN ENFOQUE CONTINUO EN EL APRENDIZAJE POR REFUERZO

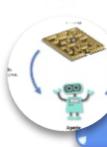
“Un Viaje a Través del Twin Delayed Deep Deterministic Policy Gradient y su Aplicación en Ambientes Desafiantes: HumanoidBulletEnv-v0”



José Javier Gutiérrez Gil
jogugil@alumni.uv.es

INTRODUCCIÓN

Resultados



Entorno



Memoria
Repetición



Explotación Vs
Exploración



Modelo



Episodio



Recompensas



Acciones

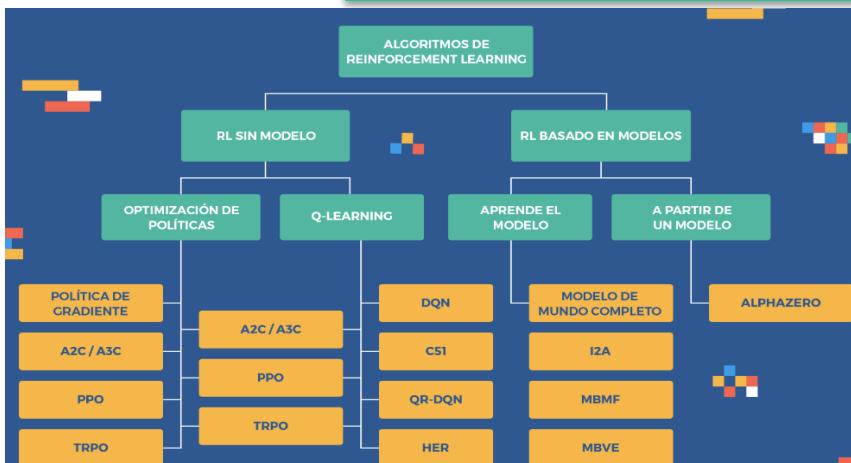


OpenAI's Gym environment
to Reinforcement Learning's atomic concepts



Introducción: RL- GYM

DDPG TD3: Gradiente Político Determinista de doble retardo



Estado

- Posición, velocidad y orientación del cuerpo del humanoide

Acciones

- **Acciones continuas** (fuerzas aplicadas a las articulaciones)

Diferencias Temporales (TD)
Vs Programación dinámica (PD):

Actualización valores Q

Optimización Política (π)

No requiere modelo completo del entorno

Entorno: PyBullet y Gym

- **HumanoidBulletEnv-v0** :Simulador físico 3D
- Aprenda a caminar de manera eficiente y sin caerse.

Recompensa:

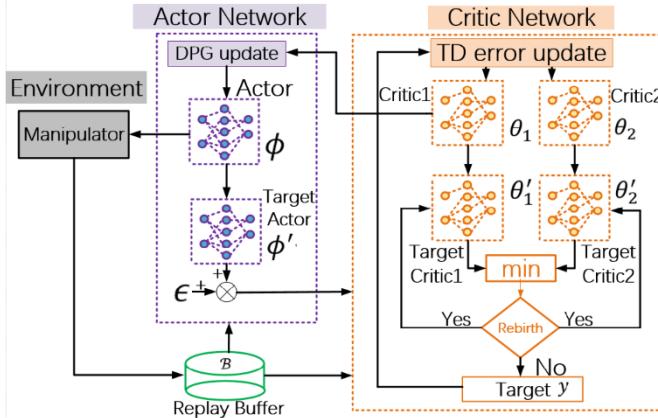
- Mantenerse de pie
- Progresión hacia delante
- Eficiencia energética



Introducción: RL- GYM

Deep Q-learning no es aplicable en el caso de acciones continuas

DDPG TD3



ReplayBuffer

(s ₁ , s' ₁ , a ₁ , r ₁)	(s ₂ , s' ₂ , a ₂ , r ₂)	(s ₃ , s' ₃ , a ₃ , r ₃)	(s ₄ , s' ₄ , a ₄ , r ₄)	...
---	---	---	---	-----

- ❑ Actor Target + Actor Model:
 - ✓ 2 Capas ocultas + 2 Relu +Tanh
- ❑ 2 Critic Target + 2 Critic Model:
 - ❑ 2 Capas Ocultas + 1 Relu

Addressing Function Approximation Error in Actor-Critic Method.
<https://arxiv.org/pdf/1802.09477.pdf>

$$R_t = \sum_{i=t}^T \gamma^{i-t} r(s_i, a_i)$$

$$J(\phi) = \mathbb{E}_{s_i \sim p_\pi, a_i \sim \pi} [R_0]$$

$$\nabla_\phi J(\phi) = \mathbb{E}_{s \sim p_\pi} [\nabla_a Q^\pi(s, a)|_{a=\pi(s)} \nabla_\phi \pi_\phi(s)]$$

$$\phi_{t+1} = \phi_t + \alpha \nabla_\phi J(\pi_\phi)|_{\phi_t}$$



TD3: Twin Delayed DDPG

2 Críticos (M Vs T) + 1 Actor (M Vs T)

Algorithm 1 Twin Delayed DDPG

```

1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi_1, \phi_2$ , empty replay buffer  $\mathcal{D}$ 
2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta$ ,  $\phi_{\text{targ},1} \leftarrow \phi_1$ ,  $\phi_{\text{targ},2} \leftarrow \phi_2$ 
3: repeat
4:   Observe state  $s$  and select action  $a = \text{clip}(\mu_\theta(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$ 
5:   Execute  $a$  in the environment
6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
8:   If  $s'$  is terminal, reset environment state.
9:   if it's time to update then
10:    for  $j$  in range(however many updates) do
11:      Randomly sample a batch of transitions,  $N = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
12:      Compute target actions
13:      
$$a'(s') = \text{clip}(\mu_{\theta_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}}), \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

14:      Compute targets
15:      
$$Q_{\phi}(r, s', d) = r + \gamma(1-d) \min_{i=1,2} Q_{\phi_{\text{targ},i}}(s', a'(s'))$$

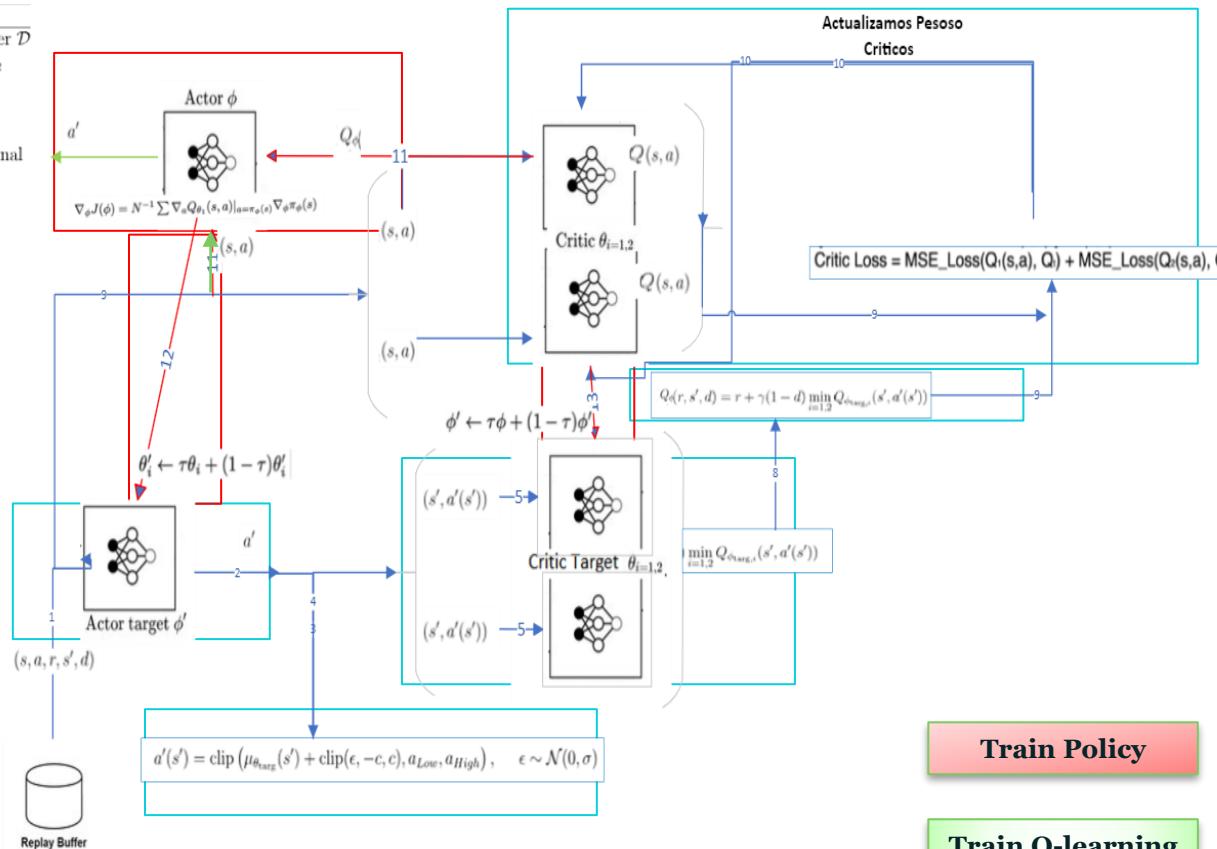
16:      Update Q-functions by one step of gradient descent using
17:      
$$\nabla_{\phi_i} \frac{1}{|N|} \sum_{(s,a,r,s',d) \in N} (Q_{\phi_i}(s, a) - Q_{\phi}(r, s', d))^2 \quad \text{for } i = 1, 2$$

18:      if  $j \bmod \text{policy.delay} = 0$  then
19:        Update policy by one step of gradient ascent using
20:        
$$\nabla_{\phi} J(\phi) = N^{-1} \sum \nabla_a Q_{\theta_1}(s, a)|_{a=\pi_\phi(s)} \nabla_\phi \pi_\phi(s)$$

21:      Update target networks with
22:      
$$\theta'_i \leftarrow \tau \theta_i + (1-\tau) \theta'_i \quad \text{for } i = 1, 2$$

23:      
$$\theta'_i \leftarrow \tau \theta_i + (1-\tau) \theta'_i$$

24:    end if
25:  end for
26: end if
27: until convergence
  
```



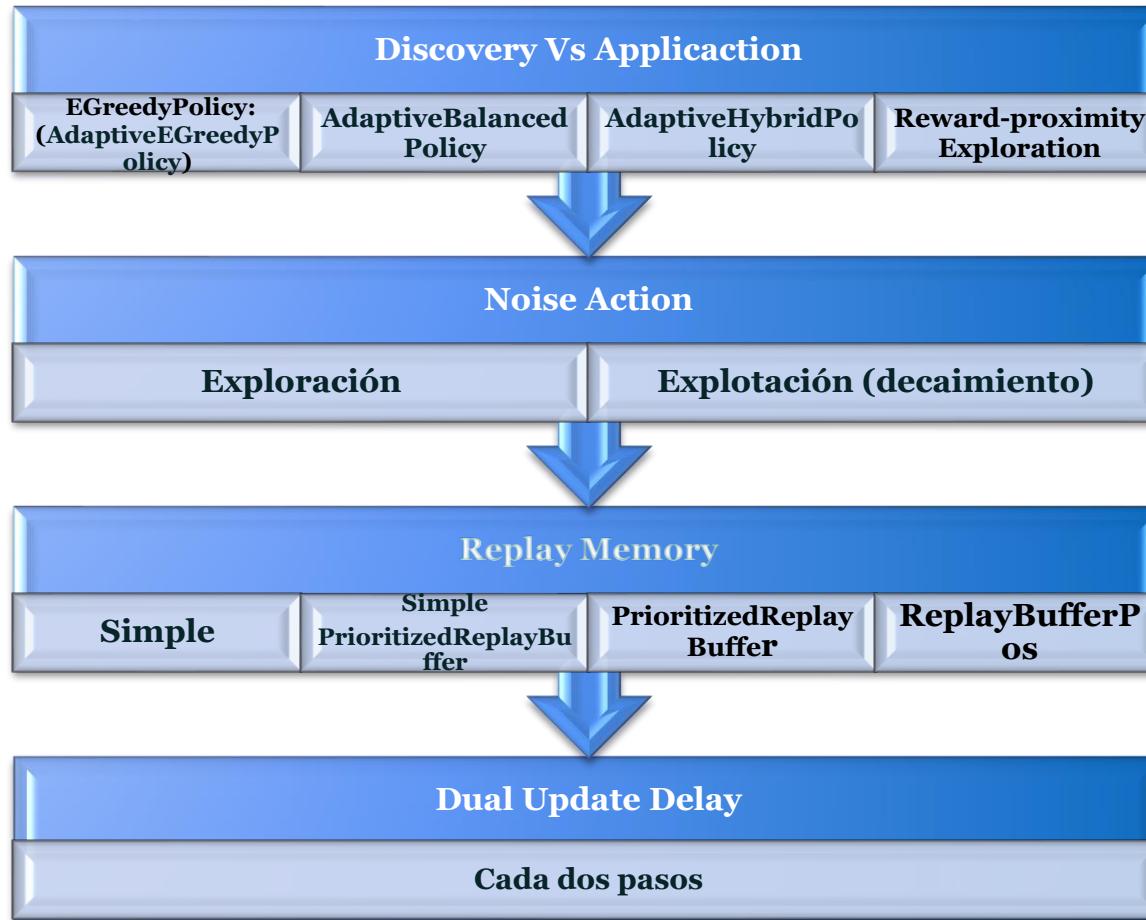
Train Policy

Train Q-learning



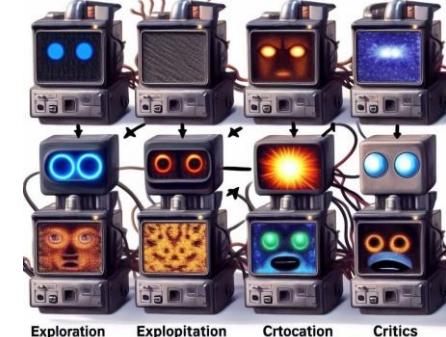
Implementación: Train

Métodos de regularización



We can might consider different methods to refine methods in order to work the model or improve it.

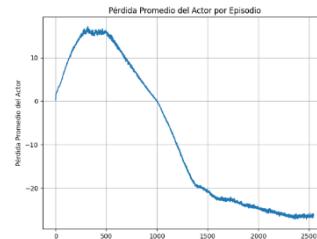
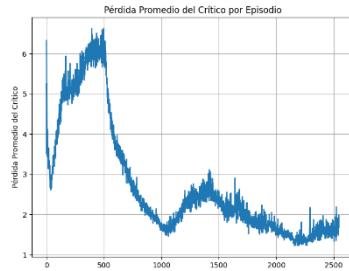
Adding noise is slowing the score, actions vs exploration at each step in that the delay is increasing.



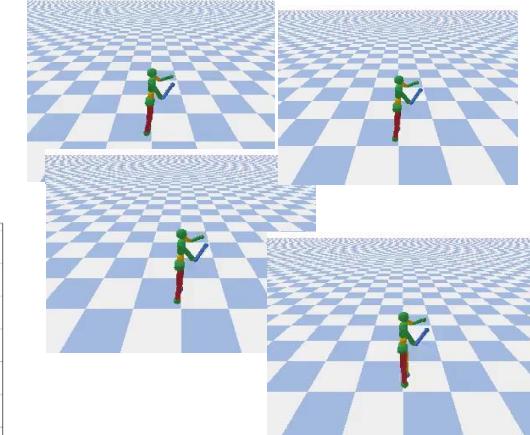
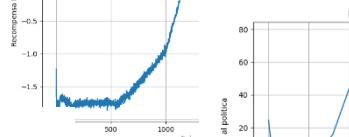
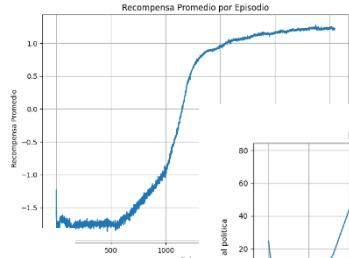


Resultados: Evaluación (I)

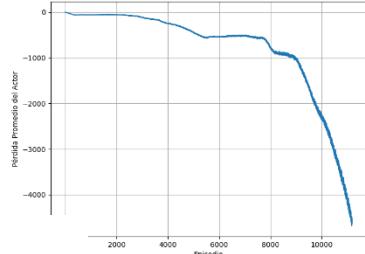
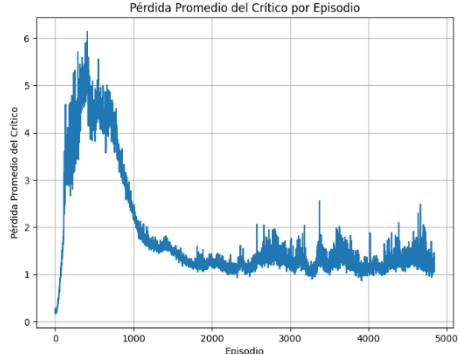
Explotación - Explotación



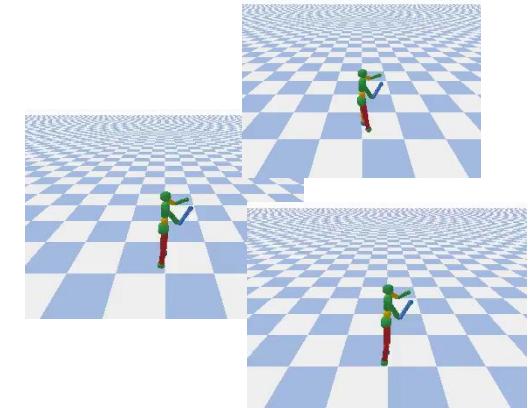
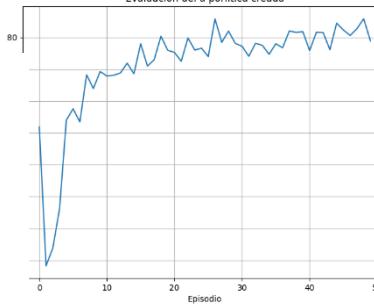
AdaptiveBalancedPolicy



Reward-proximity Exploration



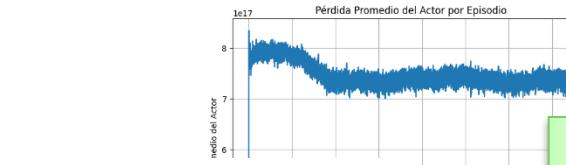
Evaluación del a política creada



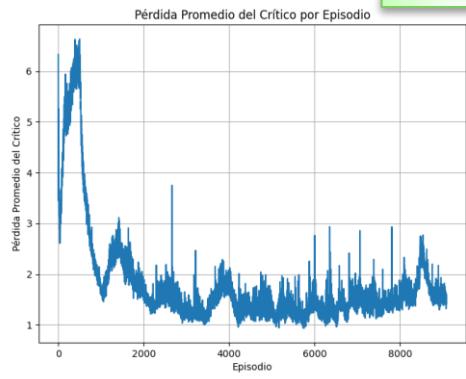
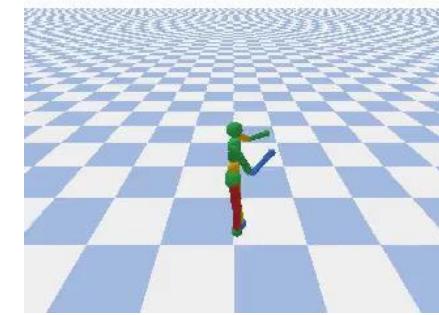
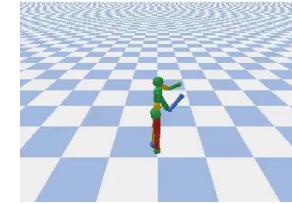
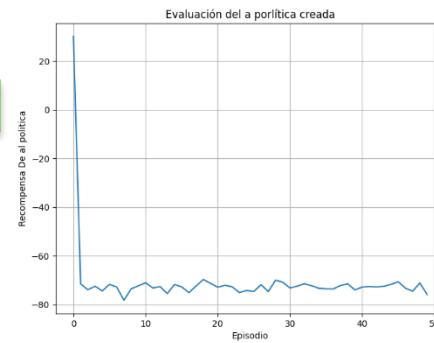
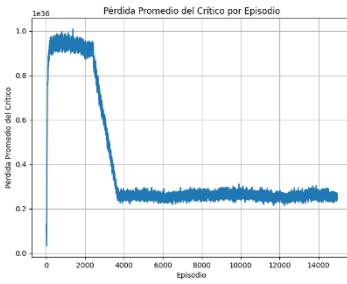


Resultados: Evaluación (II)

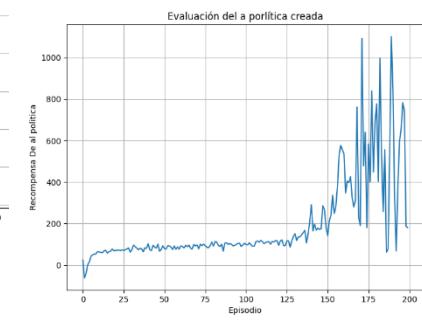
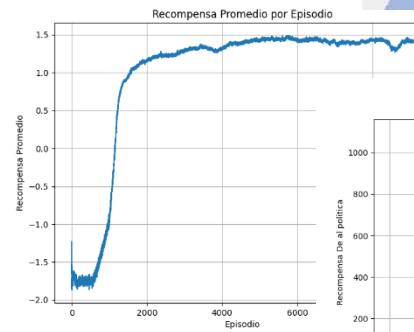
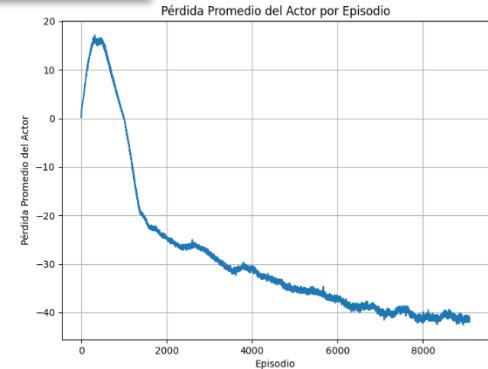
Explotación - Explotación

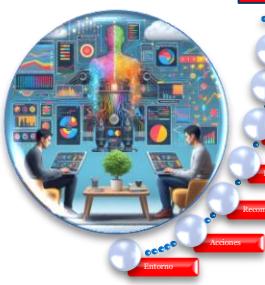


AdaptiveHybridPolicy



EGreedyPolicy

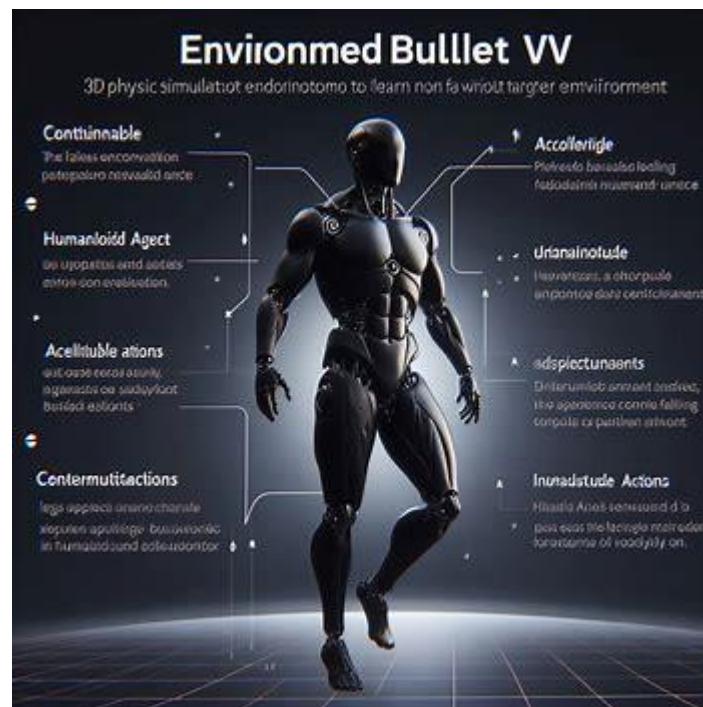




Conclusión: *Modelo RL*

Elevada Cantidad Pasos // Inestabilidad

Al evaluar la política: Se experimentan episodios en los que parece que el agente no está entrenado y su rendimiento es pobre:



- ✓ Convergencia a acciones subóptimas
- ✓ Exploración insuficiente
- ✓ Convergencia insuficiente
- ✓ Hiperparámetros incorrectos
- ✓ Problema intrínseco de la tarea



Discusión: *Modelo RL*

Elevada Cantidad Pasos // Inestabilidad

- Llegar a entrenar bien el actor - 2 críticos para que el humanoide ande decentemente es un dolor de cabeza.
- Necesidad de regulación para evitar estancamientos y sobreestimación (los críticos)
- Muy sensible a los hiperparámetros y las acciones a realizar .
- La adición de ruido $N(0,\sigma')$ mejora el algoritmo frente a búsqueda aleatoria entorno real



Aplicación: *Modelo RL*

DDPG TD3 4.670 resultados

Multiagent Copilot Approach for Shared Autonomy between Human
EEG and TD3 Deep Reinforcement Learning
<https://www.semanticscholar.org/reader/836cced5130525103c150336ae699bb36ceb86e94>

The Control Method of Twin Delayed Deep Deterministic Policy Gradient with Rebirth Mechanism to Multi-DOF Manipulator
<https://www.mdpi.com/2079-9292/10/7/870>

R. Dubey, R. Loka and A. M. Parimi, "Maintaining the Frequency of AI-based Power System Model using Twin Delayed DDPG(TD3) Implementation," 2022 2nd International Conference on Power Electronics & IoT Applications in Renewable Energy and its Control (PARC), Mathura, India, 2022, pp. 1-4, doi: 10.1109/PARC52418.2022.9726615

T. Tiong, I. Saad, K. T. K. Teo and H. b. Lago, "Deep Reinforcement Learning with Robust Deep Deterministic Policy Gradient," 2020 2nd International Conference on Electrical, Control and Instrumentation Engineering (ICECIE), Kuala Lumpur, Malaysia, 2020, pp. 1-5, doi: 10.1109/ICECIE50279.2020.9309539 •





Scott Fujimoto and Herke van Hoof and David Meger. (2028). Addressing Function Approximation Error in Actor-Critic Methods
<https://doi.org/10.48550/arXiv.1802.09477>

Kim, Myeongseop & Han, Dong-Ki & Park, Jae-Han & Kim, Jung-Su. (2020). Motion Planning of Robot Manipulators for a Smoother Path Using a Twin Delayed Deep Deterministic Policy Gradient with Hindsight Experience Replay
 Applied Sciences. 10. 575. 10.3390/app10020575.

Kim, Myeongseop & Han, Dong-Ki & Park, Jae-Han & Kim, Jung-Su. (2020). Motion Planning of Robot Manipulators for a Smoother Path Using a Twin Delayed Deep Deterministic Policy Gradient with Hindsight Experience Replay
 Applied Sciences. 10. 575. 10.3390/app10020575.

Xinrui Shen . Comparison of DDPG and TD3 Algorithms in a Walker2D Scenario. Proceedings of the 2023 International Conference on Data Science, Advanced Algorithm and Intelligent Computing (DAI 2023). Atlantis Press.D 2024.8 2024-02-13T23:00:00.000Z.P 148-155..
https://doi.org/10.2991/978-94-6463-370-2_17
 10.2991/978-94-6463-370-2_17

Pastor, J.M., Díaz, H., Armesto, L., Sala, A. Aprendizaje por refuerzo con búsqueda de políticas: simulación y aplicación a un sistema electromecánico. 7, 8 y 9 de septiembre de 2016, Madrid (pp. 710-717). DOI capítulo:
[https://doi.org/10.17979/spudc.9788497498081](https://doi.org/10.17979/spudc.9788497498081.0710)
<https://doi.org/10.17979/spudc.9788497498081>

THE-END

