

SnakeGame con Reinforcement Learning

Alan Stiven Camacho Restrepo



CONTENIDO

01 **INTRODUCCIÓN**
Reinforcement Learning y sus aplicaciones.

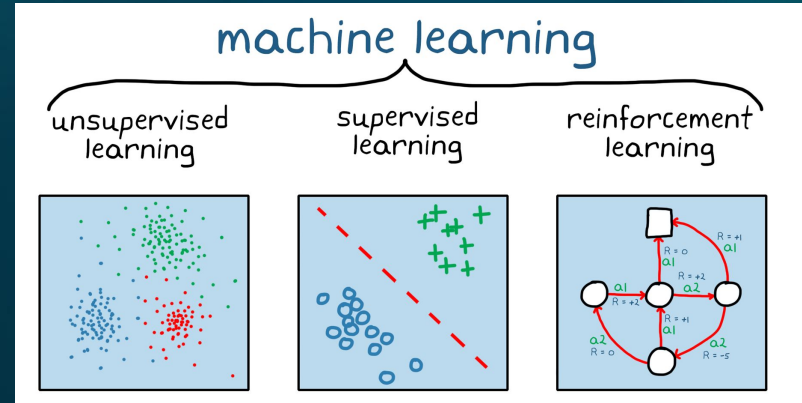
03 **METODOLOGÍA**
Construcción del entorno del juego y algoritmos de machine learning.

02 **PROBLEMA**
Planteamiento del problema, marco teórico, objetivo general y específicos.

04 **RESULTADOS**
Gráficos y videos.

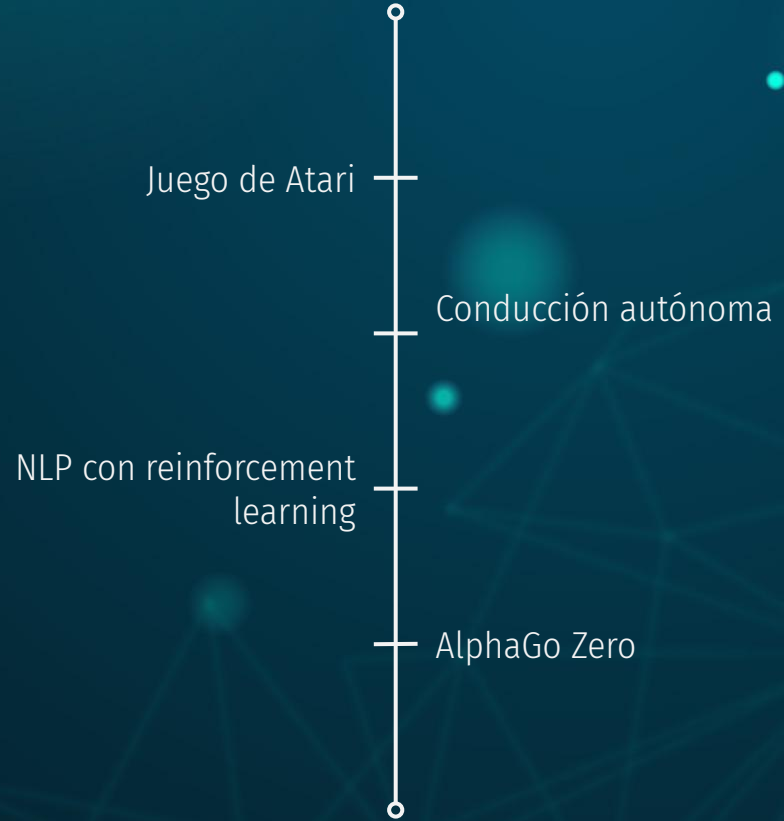
05 **CONCLUSIONES**
Perspectivas y problemas a mejorar.

¿Qué es el Reinforcement Learning?



- ❑ Aprendizaje supervisado y no-supervisado.
- ❑ ¿Hay una manera de crear un agente que aprenda a jugar por sí mismo?
- ❑ Agentes y recompensas.
- ❑ Usos más frecuentes.

APLICACIONES DEL RL



02

PLANTEAMIENTO DEL PROBLEMA

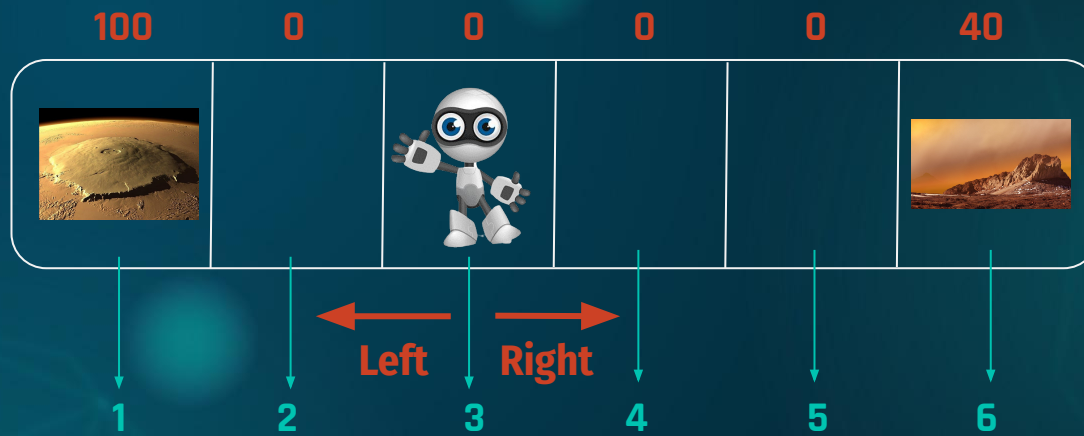
OBJETIVO GENERAL

- ❑ Construir un algoritmo de deep learning que permita que el juego de la culebra se juegue por sí misma.

OBJETIVO ESPECIFICOS

- ❑ Entorno del juego.
- ❑ Agente.
- ❑ Recompensas.
- ❑ Visualización del juego.

Reward:



- ❑ Estado.
- ❑ Acción.
- ❑ Recompensa.
- ❑ Nuevo estado.

$(s, a, R(s), s')$

Return: $R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$

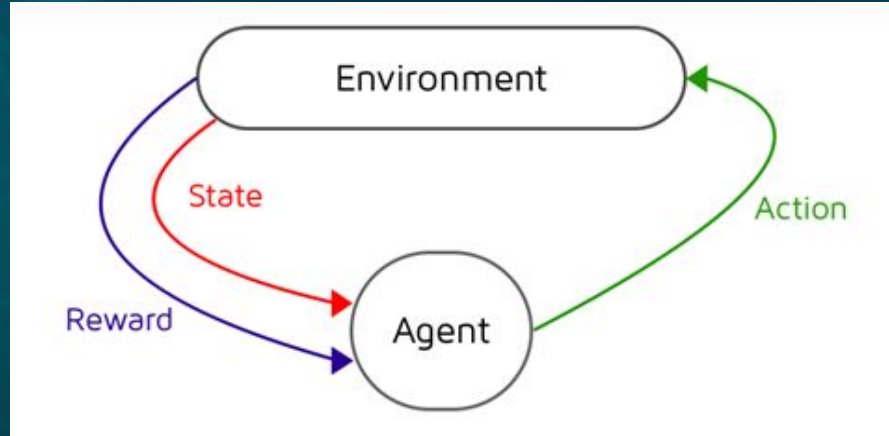
$\gamma \rightarrow$ Discount factor.
(Entre 0 y 1)

- ❑ Maximizar la función Return permite llegar a la recompensa más alta.

Policy function:

$$\pi(s) : S \rightarrow A$$

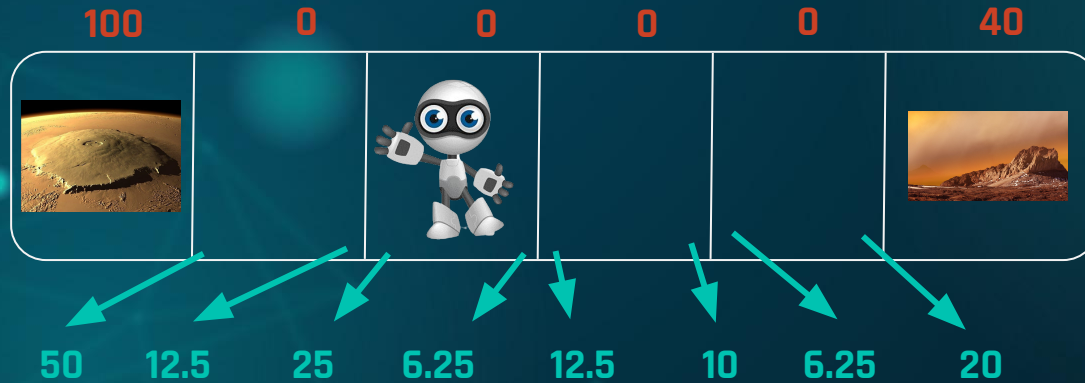
- Hallar la función que permita tomar una acción a en un estado s para que la función return se maximize.



Q-FUNCTION

$$Q(s, a) = r(s, a) + \gamma \max_a Q(s', a)$$

El mejor Return



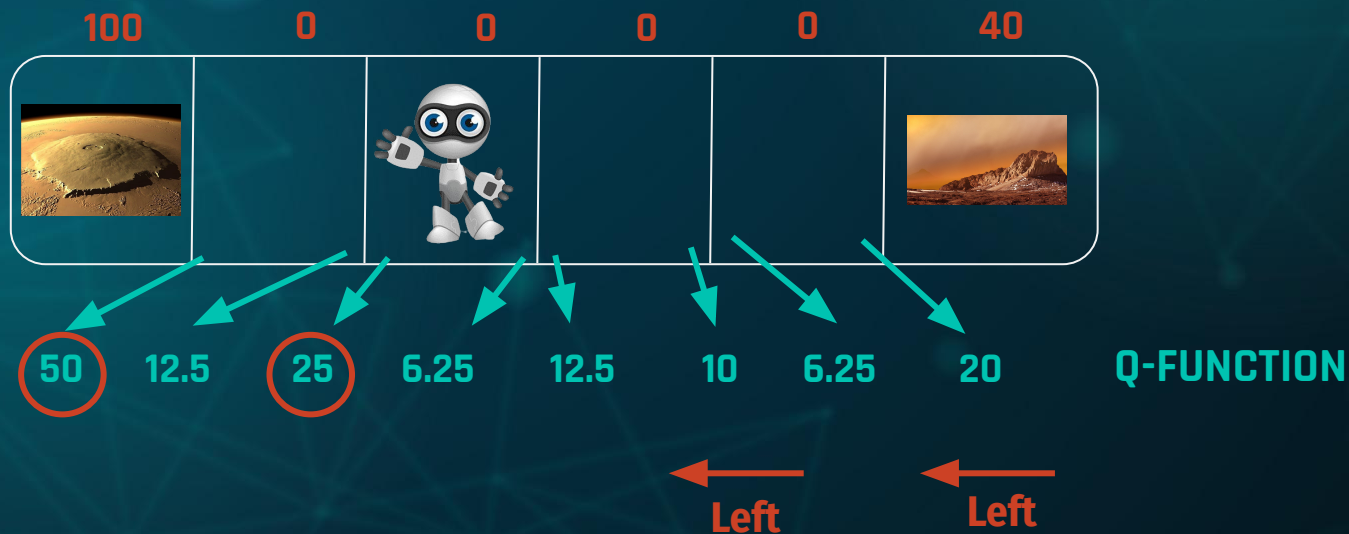
La acción más óptima en el estado s es la dada por el valor máximo de Q .

Policy
function

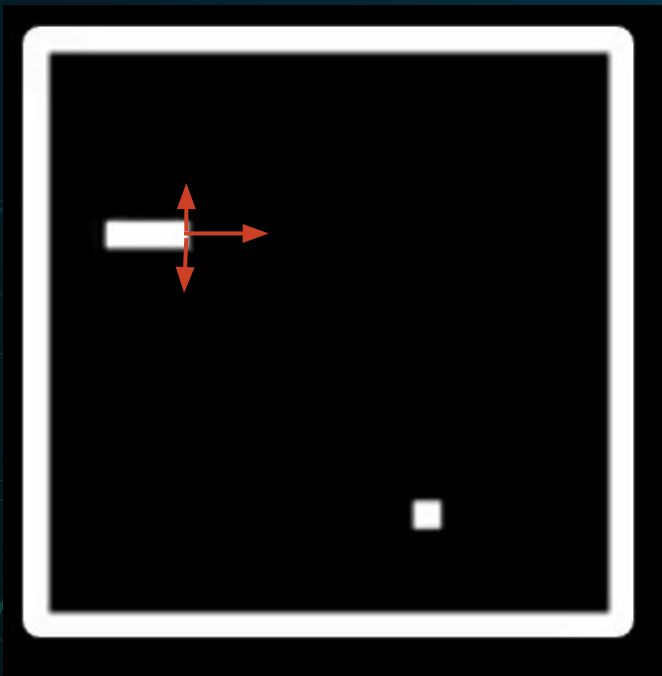
Q-function

LA META PRINCIPAL DE LA FUNCIÓN POLICY

Hallar una función que permita saber qué acción **a** se debe tomar en el estado **s** de tal forma que el agente tome el camino más óptimo a su destino.
(Proceso iterativo)



03 METODOLOGÍA

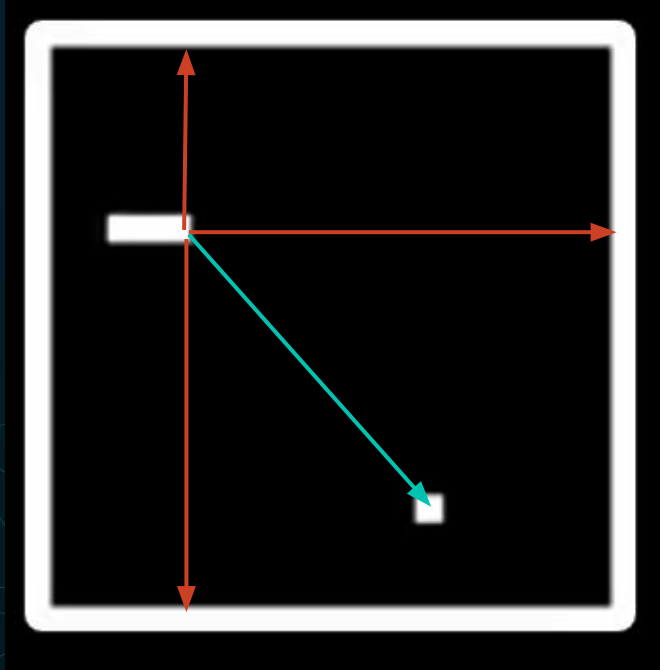


Entorno de la culebra.

Acciones: Seguir derecho = 0
 Derecha = 1
 Izquierda = 2

Recompensas: Comer alimento = + 10
 Chocar = -10
 De lo contrario = 0

ESTADO DE LA CULEBRA



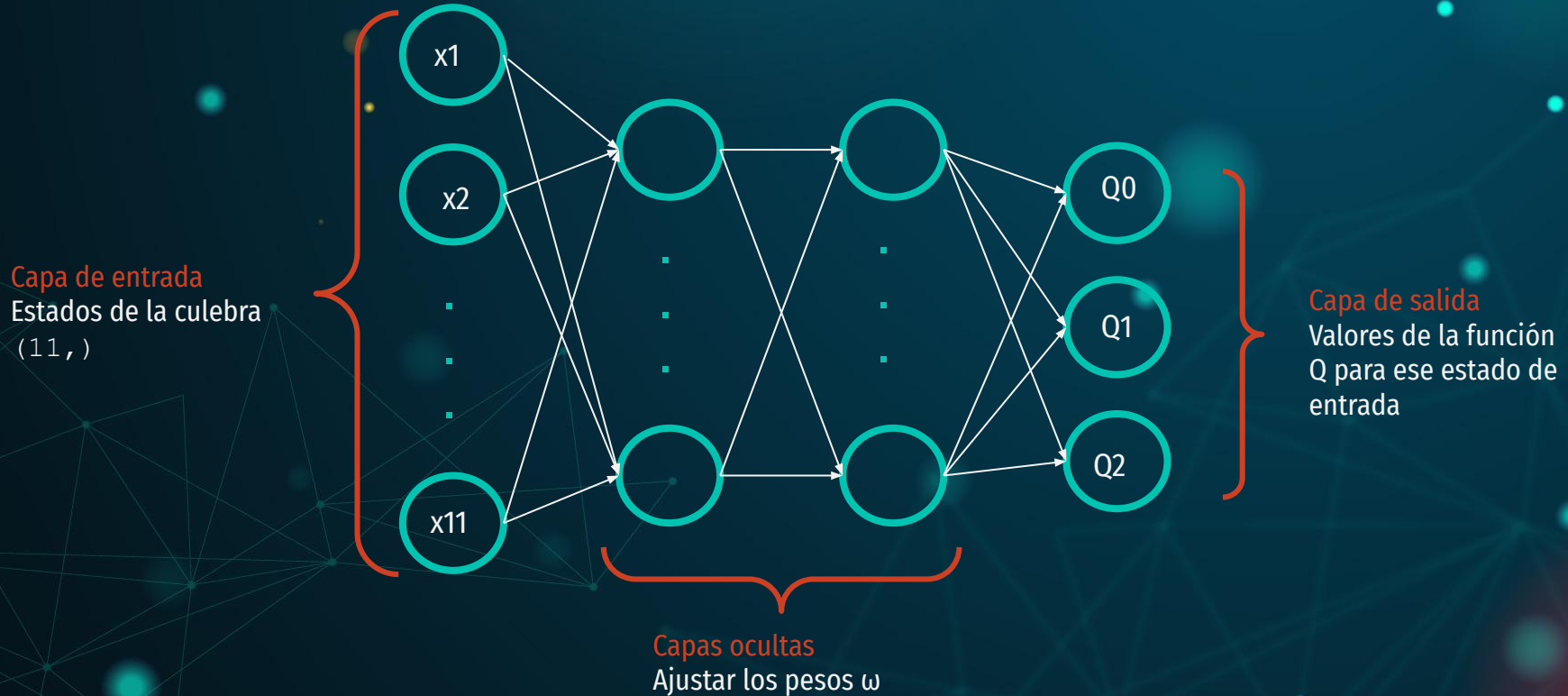
Tamaño del estado = (11,)

Peligro seguir derecho
Peligro derecha
Peligro izquierda
Dirección izquierda
Dirección derecha
Dirección arriba
Dirección abajo
Comida izquierda
Comida derecha
Comida arriba
Comida abajo

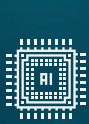
=

0
0
0
0
1
0
0
0
1
0
1

ESTRUCTURA DE LAS REDES NEURONALES



IMPLEMENTACIÓN



Uso de **dos** redes neuronales con la misma estructura.

$$\underbrace{R + \gamma \max_{a'} \hat{Q}(s', a'; w^-)}_{\text{y target}} - \underbrace{Q(s, a; w)}_{\text{Error}}$$

RED NEURONAL 1

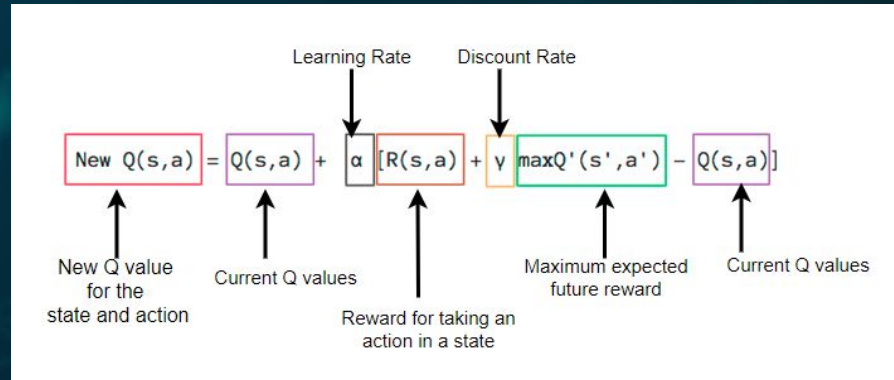
Con pesos ω'

RED NEURONAL 2

Con pesos ω

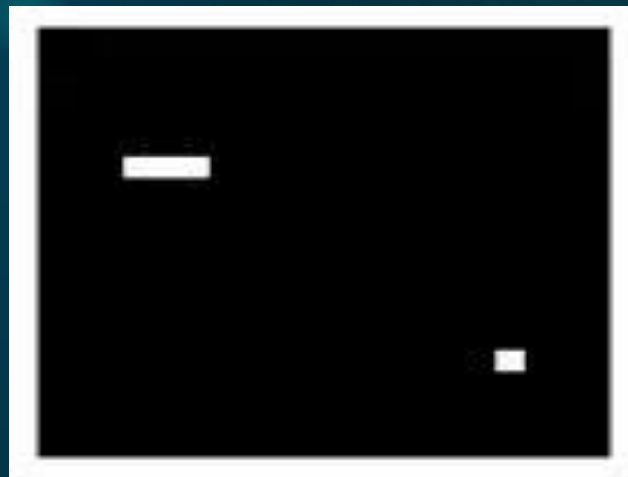
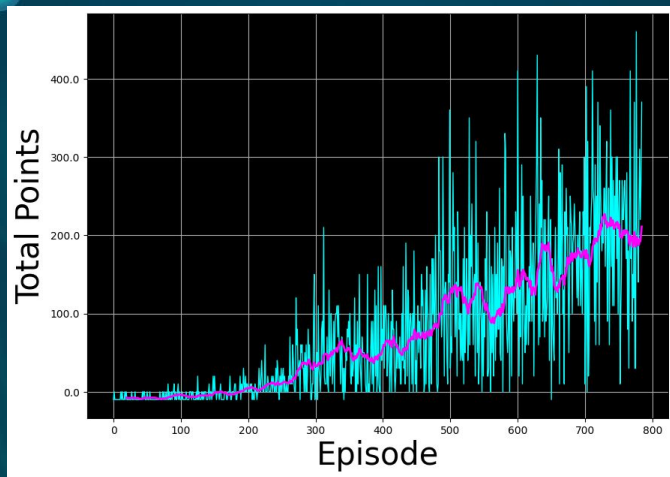
- Se inicializan los pesos aleatoriamente.
- Los pesos ω de la RN2 cambian constantemente aplicando gradiente descendiente.
- Cada C iteraciones:
 - Generar los *y targets* con la RN1.
 - Actualizar ω' con ω . (Soft update)

$$\omega' = \tau \omega + (1-\tau) \omega'$$



ALGORITMO

```
1 Initialize memory buffer  $D$  with capacity  $N$ 
2 Initialize  $Q$ -Network with random weights  $w$ 
3 Initialize target  $\hat{Q}$ -Network with weights  $w^- = w$ 
4 for episode  $i = 1$  to  $M$  do
5   Receive initial observation state  $S_1$ 
6   for  $t = 1$  to  $T$  do
7     Observe state  $S_t$  and choose action  $A_t$  using an  $\epsilon$ -greedy policy
8     Take action  $A_t$  in the environment, receive reward  $R_t$  and next state  $S_{t+1}$ 
9     Store experience tuple  $(S_t, A_t, R_t, S_{t+1})$  in memory buffer  $D$ 
10    Every  $C$  steps perform a learning update:
11    Sample random mini-batch of experience tuples  $(S_j, A_j, R_j, S_{j+1})$  from  $D$ 
12    Set  $y_j = R_j$  if episode terminates at step  $j + 1$ , otherwise set  $y_j = R_j + \gamma \max_{a'} \hat{Q}(s_{j+1}, a')$ 
13    Perform a gradient descent step on  $(y_j - Q(s_j, a_j; w))^2$  with respect to the  $Q$ -Network weights  $w$ 
14    Update the weights of the  $\hat{Q}$ -Network using a soft update
15  end
16 end
```

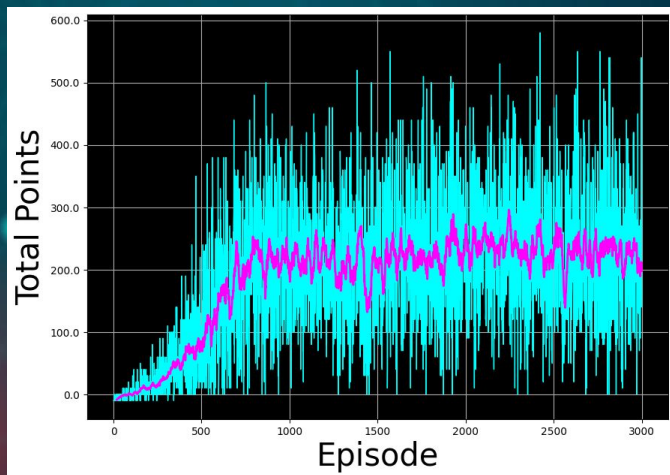


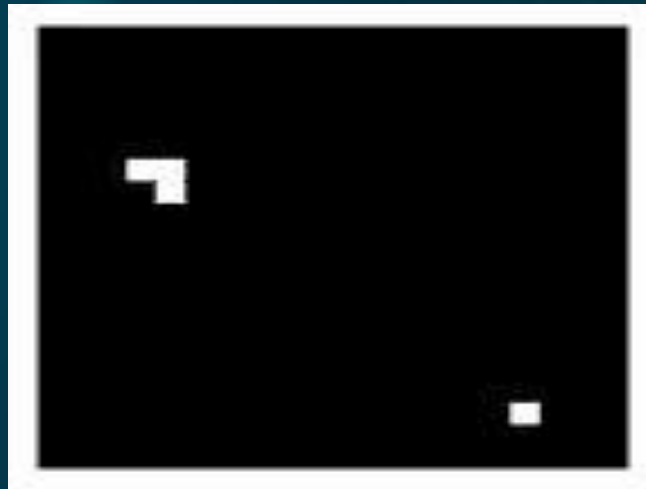
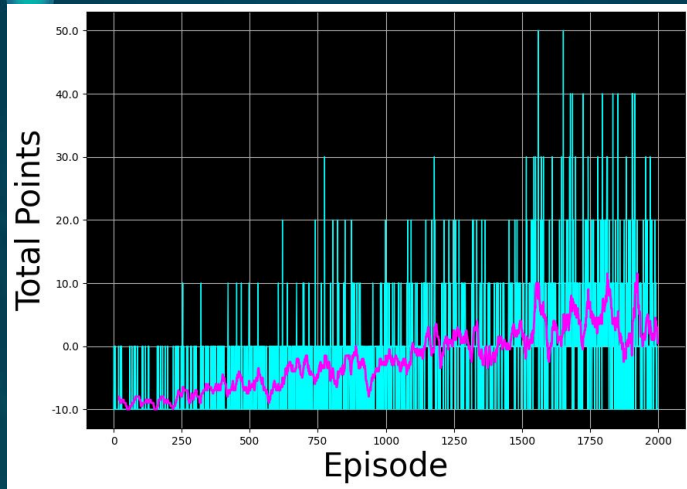
$$\alpha = 0.001$$

$$\tau = 0.001$$

$$\gamma = 0.995$$

2 Capas ocultas con 64 neuronas.



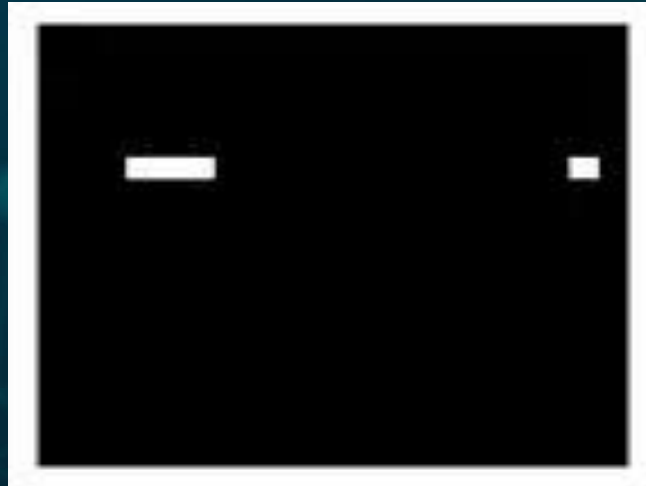
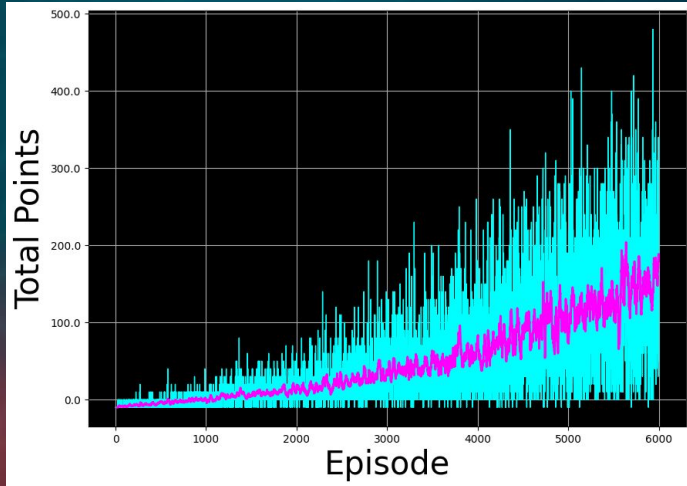


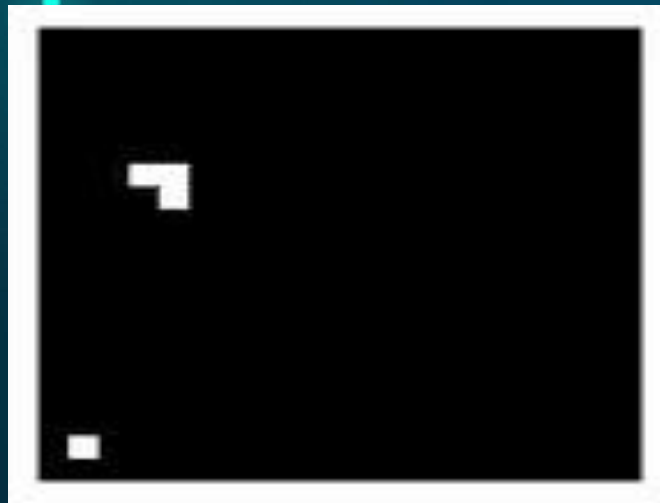
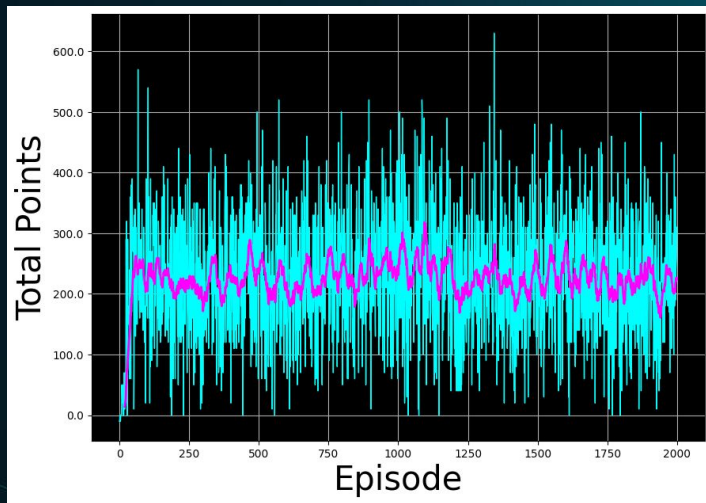
$$\alpha = 0.0001$$

$$\tau = 0.1$$

$$\gamma = 0.8$$

1 Capa oculta con 256 neuronas.



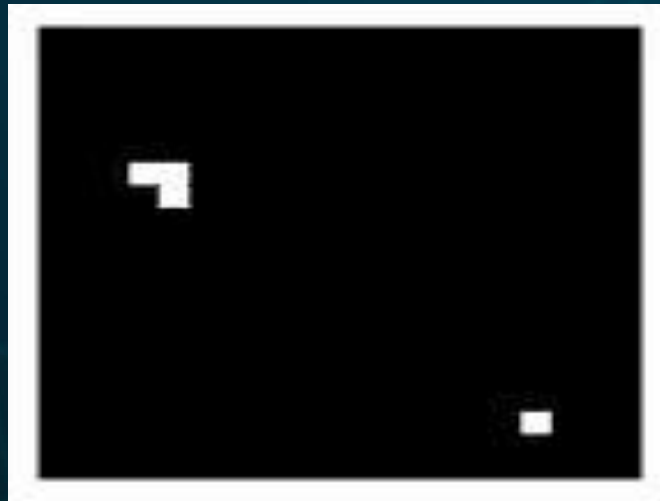
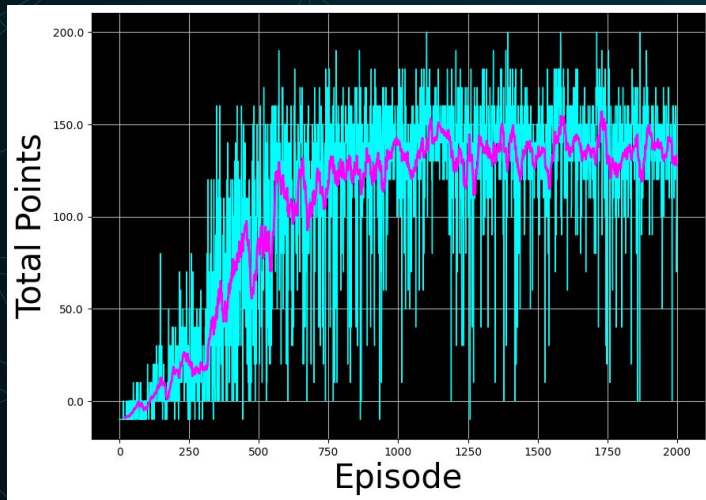


$$\alpha = 0.001$$

$$\tau = 0.01$$

$$\gamma = 0.9$$

2 Capas ocultas con 64 neuronas.



CONCLUSIONES Y PERSPECTIVAS

01	02	03	04
Se construyó una red neuronal que permitió que la culebra se jugara por sí sola.	El entorno, las recompensas, los estados y las acciones de la culebra fueron óptimos para jugar parte del juego.	Se pudo crear un entorno gráfico para la visualización del juego entrenado.	Con las redes entrenadas no se pudo completar el juego entero de la culebra.

- ❑ El reinforcement learning es un algoritmo que permite que un agente aprenda, pero es complejo saber las condiciones iniciales que debe tener el algoritmo para que haya muy buenos resultados.
- ❑ Se propone mejorar los parámetros con gráficas dependientes de las condiciones, y con las estructuras de las redes.
- ❑ Se propone usar algoritmos como los genéticos.

REFERENCIAS

- B Ravi Kiran , Ibrahim Sobh, *Deep Reinforcement Learning for Autonomous Driving: A Survey*.
- <https://neptune.ai/blog/reinforcement-learning-applications>
- <https://github.com/AleksaC/gym-snake>
- Curso de Reinforcement learning de Coursera:

<https://www.coursera.org/learn/unsupervised-learning-recommenders-reinforcement-learning?specialization=machine-learning-introduction>

THANKS!



CREDITS: This presentation template was created by **Slidesgo**, including icons by **Flaticon**, and infographics & images by **Freepik**.

Please keep this slide for attribution.