#### Reinforcement Learning: Combates de Pokémon

Antecedentes, implementación y resultados

# Alejandro Jiménez, Alejandro Gómez, Alejandro Gómez, Luis Crespo y Antonio Cabrera

Trabajo para el doble grado de Ingeniería del Software y Matemática Computacional



Asignatura de inteligencia artifial U-tad España Enero 2024

## Contents

1	Introducción				
	1.1	Contexto	3		
	1.2	Objetivos	3		
2	Explicación del juego				
	2.1	Turnos	4		
	2.2	Pokémons	4		
	2.3	Tipos	5		
		Movimientos			
3	Rei	nforcement Learning	7		
	3.1	RL	7		
	3.2	RL en Pokémon	7		

# List of Figures

2.1	Jugador eligiendo el ataque de su pokémon	4
2.2	Algunos de los pokémons.	5
2.3	Tabla de tipos	5
2.4	Un pokémon usando el movimiento hiperrayo	6

## Introducción

#### 1.1 Contexto

Desde que Deep Blue gana a Kasparov en 1997, las máquinas han ido superando a los humanos en la mayoría de juegos clásicos. Con los avances en computación y en Deep Learning la diferencia entre los grandes maestros y los ordenadores es más grande que nunca. La nueva generación de bots para jugar a juegos hace uso del Reinforcement Learning, que consiste en entrenar a una red neuronal en base a lo que aprende jugando contra ella misma.

Pokémon es la franquicia de entretenimiento más grande del mundo, con más de 300 millones de juegos vendidos. El juego consiste en capturar criaturas llamadas Pokémon y entrenarlos para que luchen contra otros Pokémon. A pesar de que el juego fue pensado para niños, Pokémon cuenta con una escena competitiva muy profesional (el torneo mundial de 2024 contará con 2.000.000 de euros en premios), sin embargo el uso de inteligencias artificiales aun no se ha extendido como en otros juegos.

#### 1.2 Objetivos

El objetivo de este trabajo es crear un bot que sea capaz de jugar a Pokémon de forma autónoma. Para ello se ha utilizado una versión simplificada del juego programada en python, en la que se eliminan ciertos elementos como el azar en la partida. El bot se ha entrenado utilizando el algoritmo de Reinforcement Q-Learning , haciendo uso de la librería de PyTorch.

## Explicación del juego

A continuación vamos a explicar de forma muy resumida las mecánicas del juego:

#### 2.1 Turnos

El juego se basa en batallas entre dos jugadores, cada uno cuenta con un equipo de 6 pokémon, de los cuales solo puede tener uno en el campo de batalla.

Las batallas se dividen en turnos, en cada turno el jugador puede realizar una acción, ya sea atacar, cambiar de pokémon. Al comienzo de cada turno ambos jugadores selecciona una acción, en el caso de que ambos decidan atacar atacará primero el pokémon con más velocidad.



Figure 2.1: Jugador eligiendo el ataque de su pokémon.

#### 2.2 Pokémons

Las batallas se realizan entre pokémons, en la actualidad existen 1025 especies de pokémons. Nosotros trabajaremos con una muestra reducida para facilitar el entrenamiento de la red neuronal.



Figure 2.2: Algunos de los pokémons.

### 2.3 Tipos

Cada pokémon tiene un tipo, existen 18 tipos diferentes, cada tipo tiene una serie de debilidades y fortalezas contra otros tipos. Además del tipo los pokémons tienen otras estadísticas como los puntos de vida, ataque, defensa o velocidad.

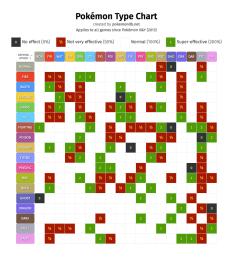


Figure 2.3: Tabla de tipos.

#### 2.4 Movimientos

Cada pokémon tiene una serie de movimientos, los cuales realizan una acción en la batalla, ya sea atacar, curar, aumentar la defensa, etc. Cada movimiento tiene un tipo, una potencia y una precisión. El tipo determina contra que pokémons es eficaz el ataque, la potencia influye en el daño del ataque y la precisión determina la probabilidad de acierto (que en nuestro caso siempre será 100%). A día de hoy existen 934 movimientos, la mayoría consisten en atacar pero hay algunas excepciones.

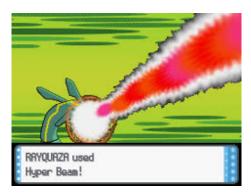


Figure 2.4: Un pokémon usando el movimiento hiperrayo.

## Reinforcement Learning

A continuación, se va a explicar brevemente en qué consiste el Reinforcement Learning (a partir de ahora RL), además de explicar cómo se aplicará este area del aprendizaje automático a nuestro trabajo.

#### 3.1 RL

El RL es un paradigma de aprendizaje automático que se centra en enseñar a un agente a tomar decisiones secuenciales para maximizar una recompensa acumulativa a lo largo del tiempo. En este enfoque de aprendizaje automático, el agente interactúa con un entorno dinámico y aprende a base de las recompensas y penalizaciones que conllevan sus decisiones.

#### 3.2 RL en Pokémon

Como en este trabajo estamos llevando a cabo batallas pokemon, el agente tiene un rango de 4 opciones a la hora de realizar una accion. Estas opciones son los 4 movimientos que puede efectuar el pokemon que está en ese momento luchando. La forma en la que se entregarán recompensas y penalizaciones al agente será en base de si estos movimientos inflingen daño (quitan puntos de vida del pokemon al que se enfrenta), o si gana el combate. De esta forma el agente aprenderá qué movimientos son los óptimos para los distintos pokemons.

El agente también posee un estado, en relación con el entorno de este trabajo. Este estado incluirá valores como, por ejemplo, el tipo o tipos de Pokémon que el agente maneja y el Pokémon al que se enfrenta, las puntuaciones de vida de los Pokémon y los movimientos disponibles de los Pokémon.

Para elegir la accion que realizará el Pokémon, el modelo que se creará tendrá que recibir como inputs, todos los valores del estado comentados previamente y devolvera un output de 4 valores correspondientes a los 4 posibles movimientos del Pokémon controlado por el agente, de los cuales se elegirá el que tenga un mayor valor.

## Deep Q-Learning

amsmath Este trabajo tambien va a estar centrado en Deep Q-Learning, un enfoque en el campo del deep learning aplicado a la toma de decisiones en entornos dinámicos como el que estamos usando. Q-Learning es un algoritmo que buscan aprender una función Q con la que evaluar la calidad de las acciones en un estado dado.

Para esto vamos a tener que implementar una red neuronal profunda para evitar utilizar tablas con todas las posibles combinaciones de acciones y estados, que en este caso serían tablas de proporciones demasiado grandes teniendo en cuenta todos los movimientos, pokemons y tipos que se han comentado anteriormente

La función Q es una función que asigna un valor para cada par estadoacción y se denota como Q(s,a), donde s respresenta el estado y a la acción. Esta función representa el valor esperado de tomar una acción en un estado específico, es decir, mide cuánto valor acumulativo (de recompensas) se espera obtener de realizar dicha acción es ese estado, y de esta forma saber cuál es la estrategia que seguir de ahí en adelante.

El Q-Learning presenta un algoritmo para entrenar al agente en el entorno propuesto, y sigue los siguientes pasos:

- 1. Inicialización
- Inicializa la función Q(s,a) de forma arbitraria para todos los pares estado-acción. Establece los parámetros del algoritmo: tasa de aprendizaje, factor de descuento y otros que se explicarán más adelante.
  - 2. Exploración/ Explotación (Realizar acciones)
- En cada frame o etapa del proceso, se elige una acción a realizar en el estado actual. Esta decisión puede ser determinista o estocástica, la primera se basa en elegir la acción que tiene el mayor valor Q conocido (explotacion), mientras que la segunda busca realizar acciones nuievas para descubrir como de efectivas o útiles son (exploración).
  - 3. Interaccion con el entorno
- Una vez la acción ha sido realizada, se observa la recompensa o penalización recibida y el próximo estado.

- 4. Actualización de la función Q
- Se usa la ecuación de Bellman para actualizar Q(s,a):

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot \left[ r + \gamma \cdot \max_{a'} Q(s',a') - Q(s,a) \right]$$
(4.1)

Donde: - (alpha) es la tasa de aprendizaje que controla la magnitud de la actualización. - r es la recompensa obtenida - gamma como el factor de descuento que penaliza las recompensas a largo plazo. -  $\max_a'Q(s',a')$  representa el valor máximo esperado en el proximo estado s'

- 5. Repetición
- Se realizan los pasos del 2 al 4 durante un número de frames o iteraciones.
- 6. Convergencia
- La función Q converge a la función de valor óptimo a medida que el agente explora y aprende mas sobre el entorno.
  - 7. Política óptima
- La política óptima se obtiene a partir de la función Q resultante, la acción óptima es aquella que maximiza Q(s,a).

Este algoritmo es un proceso iterativo en el que el agente interactúa con el entorno y cambia sus estimaciones de Q, de esta forma creando una política o estrategia para maximizar las recompensas. Un equilibrio entre la exploración y la explotación es de gran importancia para asegurar que el agente conoce las mejores acciones a realizar sin atascarse en comportamientos subóptimos.