

Alpha Zero.

Edwin Alexander Enciso, Luis David Restrepo Cadavid,
Sebastián Lopera Osorio
david95rc@utp.edu.co, edwinenciso@utp.edu.co
sebastian.lopera@utp.edu.co

Resumen- El futuro ha llegado de golpe. Hace veinte años muchos dudaban que una máquina fuese capaz de derrotar a los humanos en una partida de ajedrez. La victoria de Deep Blue, el ordenador creado por IBM, en 1997 sobre el campeón del mundo Gary Kasparov supuso un cambio de paradigma. Desde ese momento, la historia del deporte-ciencia cambió para siempre. Sin embargo, hace una semana llegó el holocausto absoluto. Un programa de inteligencia artificial destrozó de forma humillante al campeón de los ordenadores ajedrecistas, cambiando para siempre la disciplina y abriendo una nueva era a nivel científico.

Palabras clave- ajedrez, Deep Blue, inteligencia artificial.

Abstract— The future has come suddenly. Twenty years ago many doubted that a machine was capable of defeating humans in a game of chess. The victory of Deep Blue, the computer created by IBM, in 1997 over world champion Gary Kasparov was a paradigm shift. From that moment on, the history of sports science changed forever. However, a week ago the absolute holocaust arrived. A program of artificial intelligence humiliatingly destroyed the champion of chess computers, forever changing discipline and opening a new era at the scientific level.

Keywords- Chess, Deep Blue, Artificial intelligence.

I. INTRODUCCION

Han pasado veinte años desde aquello. Hoy, uno de los campeones más implacables de la IA es AlphaZero, un programa de ordenador desarrollado por la compañía DeepMind y que nace de AlphaGo, que fue capaz de derrotar al mejor jugador de go del mundo. AlphaZero es capaz de alcanzar un nivel de maestría sobrehumana en el ajedrez, el go y el shogi, y además lo hace aprendiendo. Puede aprender desde cero, tan solo conociendo las reglas, en solo 24 horas, sencillamente por el hecho de jugar contra sí mismo. Ya no necesita que nadie programe reglas que debe seguir, o hacer infinidad de cálculos antes de cada jugada.

Hace un año, DeepMind publicó un borrador en el que demostraba que AlphaZero, era capaz de derrotar a los programas más avanzados para jugar a los tres juegos comentados: Stockfish, Elmo y el propio AlphaGo Zero. Lo hizo usando circuitos diseñados para permitir el aprendizaje maquinal y basados en redes neurales. Todos ellos le permitían a esta IA no solo alcanzar un rendimiento extraordinario, sino también acomodarse a una gama más amplia de reglas. Estos logros acaban de publicarse en la revista Science.



Resultados de AlphaZero jugando contra Stockfish, Elmo y AlphaGo Zero – DeepMind.

Nuestros resultados demuestran que un algoritmo de aprendizaje por refuerzo y de propósito general puede aprender, partiendo de cero –sin necesidad de añadir conocimientos o datos previos, proporcionados por humanos, solo conociendo las reglas– y alcanzar un rendimiento sobrehumano en varios juegos de gran complejidad», han escrito los autores de la investigación, dirigida por David Silver. Por tanto, AlphaZero representa un importante paso adelante en la tarea de crear una avanzada inteligencia artificial capaz de dominar juegos más complejos por su cuenta.

El objetivo de DeepMind es construir sistemas que puedan solucionar algunos de los problemas más complejos del mundo y crear un programa que puede enseñarse a sí mismo cómo dominar el ajedrez, el shogi y el go desde cero es un importante primer paso en ese camino», ha dicho en un comunicado Demis Hassabis, director y cofundador de DeepMind.

II. CONTENIDO

En la teoría de juegos, más que razonar sobre juegos específicos, a los matemáticos les gusta razonar sobre una clase especial de juegos: juegos de dos jugadores por turnos con información perfecta. En estos juegos, ambos jugadores saben todo lo relevante sobre el estado del juego en un momento dado. Además, no hay aleatoriedad o incertidumbre sobre cómo los movimientos afectan el juego; hacer un movimiento dado siempre dará como resultado el mismo estado final del juego, uno que ambos jugadores conocen con total certeza. Llamamos juegos como estos juegos clásicos.

Estos son ejemplos de juegos clásicos:

Tic-Tac-Toe
Ajedrez
Ir
Gomoku (un juego de 5 en fila en un tablero de go 19 por 19)
Mancala
Mientras que estos juegos no son:

Póker (y la mayoría de los otros juegos de cartas)
 Piedra Papel tijeras
 El dilema del prisionero
 Videojuegos como Starcraft
 Cuando los juegos tienen elementos aleatorios y estados ocultos, es mucho más difícil diseñar sistemas de IA para jugarlos, aunque se han desarrollado potentes póquer y Starcraft AI. Por lo tanto, el algoritmo AlphaZero se limita a resolver solo juegos clásicos.[4]

Sin embargo con unas pequeñas mejoras es capaz de adaptarse a cualquier tipo de juego. En marzo de 2016, AlphaGo de Deepmind venció 18 veces al campeón mundial de Go player Lee Sedol 4-1 en una serie vista por más de 200 millones de personas. Una máquina había aprendido una estrategia sobrehumana para jugar Go, una hazaña que antes se creía imposible, o al menos, al menos a una década de lograrse.[5]



Imagen 2[5]

Esto en sí mismo, fue un logro notable. Sin embargo, el 18 de octubre de 2017, DeepMind dio un salto gigante más allá. El artículo 'Dominando el juego de Go sin conocimiento humano' reveló una nueva variante del algoritmo, AlphaGo Zero, que había derrotado a AlphaGo 100-0. Increíblemente, lo hizo aprendiendo únicamente a través del juego propio, comenzando la 'tabula rasa' (estado en blanco) y encontrando gradualmente estrategias que superaran las encarnaciones anteriores de sí mismo. Ya no se requería una base de datos de juegos expertos humanos para construir una IA superhumana.[5]

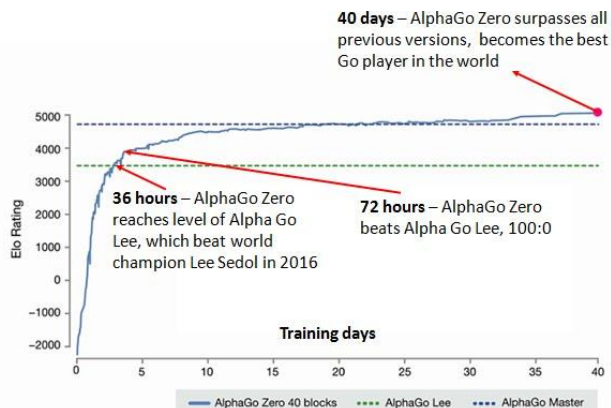


Imagen 3 [5]

Apenas 48 días después, el 5 de diciembre de 2017, DeepMind publicó otro artículo 'Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm' que mostraba cómo AlphaGo Zero podía adaptarse para vencer a los programas campeones del mundo StockFish y Elmo en ajedrez y shogi. Todo el proceso de aprendizaje, desde mostrar los juegos por primera vez, hasta convertirse en el mejor programa de computadora del mundo, había llevado menos de 24 horas.

Con esto nació AlphaZero - el algoritmo general para ser bueno en algo, rápidamente, sin ningún conocimiento previo de la estrategia de expertos humanos.[5]

Cómo construir tu propio AlphaZero

En primer lugar, consulte la hoja de trucos de AlphaGo Zero para obtener un alto nivel de comprensión de cómo funciona AlphaGo Zero. Vale la pena referirnos a eso a medida que recorremos cada parte del código. También hay un gran artículo aquí que explica cómo funciona AlphaZero con más detalle.[5]

El código

Clona este repositorio Git, que contiene el código que voy a referenciar.

Para iniciar el proceso de aprendizaje, ejecute los dos paneles superiores del cuaderno.ipynb Jupyter. Una vez que haya construido suficientes posiciones de juego para llenar su memoria, la red neuronal comenzará a entrenar. A través del auto-juego y entrenamiento adicional, gradualmente mejorará en la predicción del valor del juego y los siguientes movimientos desde cualquier posición, resultando en una mejor toma de decisiones y un juego más inteligente en general. [5]

Ahora veremos el código con más detalle, y mostraremos algunos resultados que demuestran que la IA se hace más fuerte con el tiempo.

N.B - Este es mi propio entendimiento de cómo funciona AlphaZero basado en la información disponible en los documentos mencionados anteriormente. Si alguno de los siguientes puntos es incorrecto, ¡disculpe y me esforzaré por corregirlo! [5]

Conectar4

El juego que nuestro algoritmo aprenderá a jugar es Connect4 (o Four In A Row). No tan complejo como el Go... pero todavía hay 4.531.985.219.092 posiciones de juego en total.



Imagen 4 [5]

Las reglas del juego son sencillas. Los jugadores se turnan para introducir una pieza de su color en la parte superior de cualquier columna disponible. El primer jugador que consiga cuatro de su color seguidos - cada uno vertical, horizontal o diagonalmente - gana. Si toda la cuadrícula se llena sin que se cree una cuadrícula de cuatro en una fila, el juego se empata.

He aquí un resumen de los archivos clave que componen la base de código:

game.py

Este archivo contiene las reglas del juego para Connect4.

A cada cuadrado se le asigna un número del 0 al 41, como sigue:

0	1	2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31	32	33	34
35	36	37	38	39	40	41

Imagen 5 [5]

El archivo `game.py` da la lógica para pasar de un estado de juego a otro, dada una acción elegida. Por ejemplo, dado el tablero vacío y la acción 38, el método `takeAction` devuelve un nuevo estado de juego, con la ficha del jugador titular en la parte inferior de la columna central.

Puede reemplazar el archivo `game.py` por cualquier archivo de juego que se ajuste a la misma API y el algoritmo, en principio, aprenderá la estrategia a través del auto-juego, basado en las reglas que usted le ha dado.[5]

run.ipynb

Contiene el código que inicia el proceso de aprendizaje. Carga las reglas del juego y luego itera a través del bucle principal del algoritmo, que consta de tres etapas:

1. **Auto-juego**
2. **Reentrenamiento de la red neuronal**
3. **Evaluación de la red neuronal**

Hay dos agentes implicados en este bucle, el `mejor_jugador` y el `current_player`.

El `best_player` contiene la red neuronal de mejor rendimiento y se utiliza para generar las memorias de auto-juego. El `current_player` entonces vuelve a entrenar su red neuronal en estas memorias y es entonces lanzado contra el `best_player`. Si gana, la red neuronal dentro del `best_player` es cambiada por la red neuronal dentro del `current_player`, y el bucle comienza de nuevo.[5]

agent.py

Contiene la clase `Agent` (un jugador en el juego). Cada jugador se inicia con su propia red neuronal y el árbol de búsqueda de Monte Carlo.

El método de simulación ejecuta el proceso de búsqueda en el árbol de Monte Carlo. Específicamente, el agente se mueve a un nodo de hoja del árbol, evalúa el nodo con su red neural y luego rellena el valor del nodo a través del árbol.

El método de acción repite la simulación varias veces para comprender cuál es el movimiento más favorable desde la posición actual. A continuación, devuelve la acción elegida al juego, para realizar la jugada.

El método de repetición vuelve a entrenar la red neuronal, utilizando memorias de juegos anteriores.[5]

Model.py

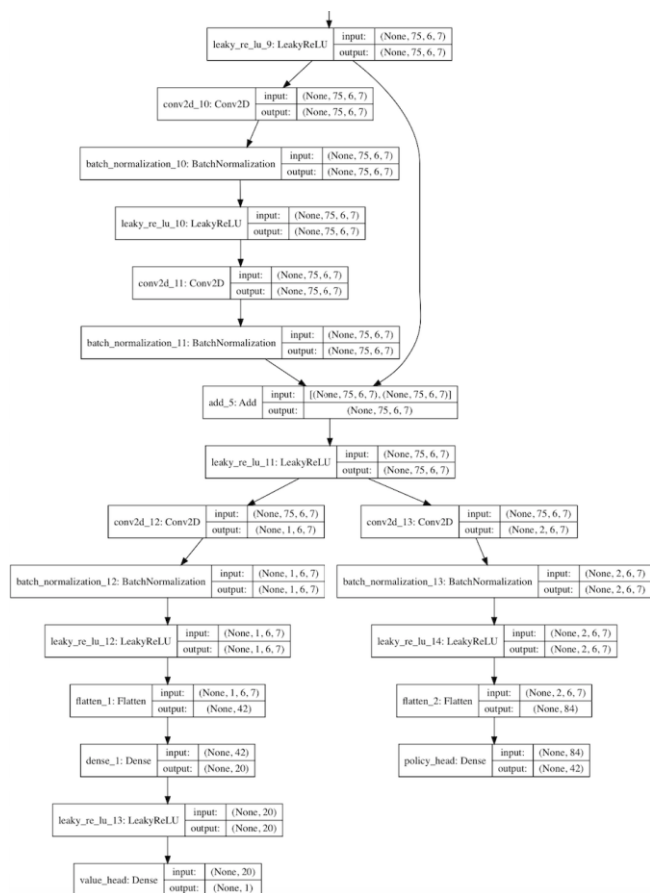


Imagen 6 [5]

Este archivo contiene la clase Residual_CNN, que define cómo construir una instancia de la red neuronal.

Utiliza una versión condensada de la arquitectura de la red neuronal en el documento AlphaGoZero, es decir, una capa convolucional, seguida de muchas capas residuales, que luego se dividen en un valor y un jefe de política.

La profundidad y el número de filtros convolucionales se pueden especificar en el archivo de configuración.

La biblioteca Keras se utiliza para construir la red, con un backend de Tensorflow.

Para ver filtros convolucionales individuales y capas densamente conectadas en la red neural, haga lo siguiente dentro del cuaderno de run.ipynb:

```
current_player.model.viewLayers()
```

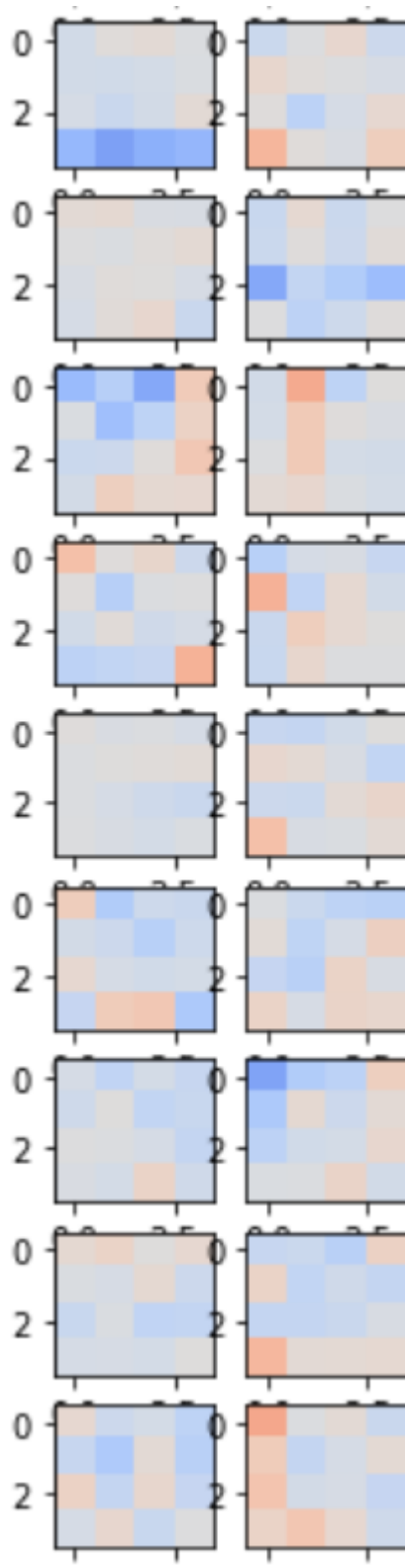


Imagen 7 [5]

MCTS.py

Contiene las clases Node, Edge y MCTS, que constituyen un árbol de búsqueda de Monte Carlo.

La clase MCTS contiene los métodos moveToLeaf y backFill mencionados anteriormente, e instancias de la clase Edge almacenan las estadísticas sobre cada movimiento potencial.

4. <https://nikcheerla.github.io/deeplearningschool/2018/01/01/AlphaZero-Explained/>
5. <https://medium.com/applied-data-science/how-to-build-your-own-alphazero-ai-using-python-and-keras-7f664945c188>