

Reinforcement Learning aplicado en la Solución Automática de un Laberinto

Alejos Arrieta, Fredo Antonino
20190390D
Universidad Nacional de Ingeniería
Lima, Perú
fredo.alejos.a@uni.pe

Diaz Ramos, Noe Elvis
20181362A
Universidad Nacional de Ingeniería
Lima, Perú
noe.diaz.r@uni.pe

I. ENUNCIADO DEL CASO

Diseño de un laberinto resuelto automáticamente mediante redes neuronales y programado en Python

II. FUNCIONAMIENTO DEL SOFTWARE

Es una implementación de un juego de laberinto en Pygame. Las clases Player y Wall se utilizan para representar al jugador y a las paredes del laberinto, respectivamente. La función arraymax devuelve una lista de índices de elementos máximos en una lista dada. El nivel del laberinto se representa como una lista de cadenas, y se utiliza para crear una lista de paredes y un rectángulo para el final del laberinto.

Después de definir el laberinto, el código define los estados posibles y las acciones que puede tomar el jugador. Los estados se representan como coordenadas (x, y) y las acciones como cadenas "UP", "DOWN", "LEFT", y "RIGHT".

La función de recompensa, reward, toma un estado y una acción como entrada y devuelve una recompensa. Si el estado es el mismo que la posición del final del laberinto, se devuelve una recompensa de 100. De lo contrario, se devuelve una recompensa de 0.

III. LIMITACIONES DEL SOFTWARE

El laberinto se define como una lista de cadenas, lo que puede ser poco práctico si queremos modificar el laberinto de alguna manera. Puede ser más fácil representar el laberinto como una matriz de celdas o una lista de objetos Wall y End en lugar de una lista de cadenas.

El código no maneja bien la colisión entre el jugador y el final del laberinto. En lugar de devolver la recompensa de 100 cuando el jugador llegue al final, el código debería detener el juego o volver al menú principal o algo similar.

El código no tiene ningún mecanismo para elegir una acción para el jugador. Esto podría ser una limitación importante si quieres implementar algún tipo de inteligencia artificial o aprendizaje automático para que el jugador tome decisiones.

El código no tiene un bucle de juego principal.

IV. MARCO TEÓRICO

A. Aprendizaje por refuerzo

El aprendizaje por refuerzo es una rama del machine learning en la cual la máquina guía su propio aprendizaje a través



REINFORCEMENT LEARNING MODEL

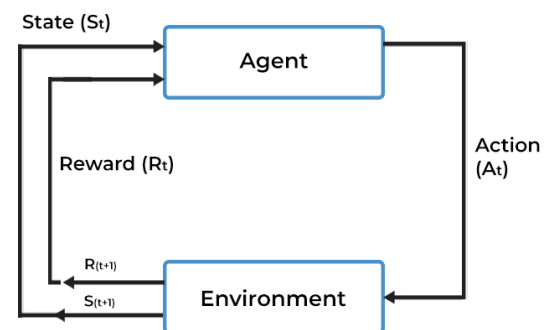


Fig. 1. Modelo básico del aprendizaje por refuerzo. Fuente: Adaptado de [2]

de recompensas y castigos. Es decir, consiste en un sistema de instrucción autónomo cuyo camino es indicado según sus aciertos y errores.

Consta de un aprendizaje empírico, por lo que el agente informático está en constante búsqueda de aquellas decisiones que le premien de algún modo, a la par que evita aquellos caminos que, por experiencia propia, son penalizados.

También, se puede decir que el aprendizaje reforzado es un concepto similar al que utilizan los seres vivos. Esto es, las máquinas aprenden qué decisiones tomar de acuerdo a la situación en la que se encuentren. Además, son capaces de desarrollar estrategias con una visión a largo plazo. [1].

B. Q-Learning

Se trata de una técnica de aprendizaje por refuerzo. Consiste en que el sistema aprenda una serie de normas que le diga al agente que acción bajo que circunstancias. Entre sus características principales se encuentra en que no requiere un modelo de entorno y puede manejar problemas con transiciones estocásticas y recompensas sin requerir adaptaciones.

1) *Funcionamiento*: El algoritmo del Q-Learning funciona con un conjunto de datos llamada matriz Q, que se actualiza

progresivamente mediante la siguiente función.

$$Q_t^{new} = (1 - \alpha)Q_t^{old} + \alpha(r_t + \gamma \max_a(Q_{t+1})) \quad (1)$$

Donde, α representa la tasa de aprendizaje, γ representa el factor de descuento, el cual es encargado de evaluar las recompensas recibidas anteriormente con un valor mayor que las recibidas posteriormente. Finalmente $Q(t)$ es el valor obtenido de Q según el estado y acción aplicado en el tiempo actual t . Es importante que el significado cualitativo de $\max_a(Q_{t+1})$ es un estimado de un valor óptimo en el futuro.[3]

V. DIAGRAMA DE CLASES

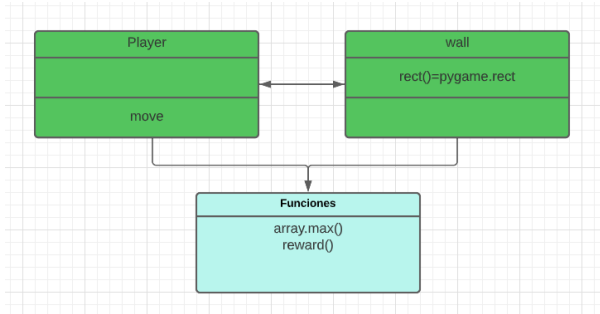


Fig. 2. Diagrama de Clases

VI. ALGORITMO DE LAS FUNCIONALIDADES DEL SOFTWARE

La clase Player representa al jugador y tiene un método move que se utiliza para mover al jugador en una sola dirección. La clase Wall representa una pared y solo tiene un atributo, rect, que es un objeto pygame.Rect que se utiliza para representar la posición y el tamaño de la pared. La función arraymax toma una lista como entrada y devuelve una lista de índices de elementos máximos en la lista. La lista de cadenas level contiene el diseño del laberinto. La función reward toma un estado y una acción como entrada y devuelve una recompensa. El código inicializa Pygame y crea una ventana para dibujar el laberinto y el jugador. El código crea una lista de paredes y un rectángulo para el final del laberinto a partir de la lista de cadenas level. Los estados y las acciones se definen como variables globales.

VII. DESCRIPCIÓN DE LAS FUNCIONES

Dentro de las funciones que se emplearon se encuentran las siguientes:

- **Dentro de el objeto Player** En el objeto player se encuentran las siguientes funciones:
 - **move(self, dx, dy):** Esta función se encarga de trasladar la ubicación del jugador, mediante un desplazamiento indicado por dx y dy, esta hace uso de la función **move_single_axis(self, dx, dy)**.
 - **move_single_axis(self, dx, dy):** Esta función se encarga de trasladar el jugador en una sola dirección,

garantizando que no habra solapamiento entre el jugador y los muros del juego.

- **arraymax(array)** Esta función se encarga de obtener los índices de todos los elementos con el máximo valor dentro del arreglo **array**. Esto se emplea para la elección de dirección del jugador durante el aprendizaje.

VIII. DESCRIPCIÓN DE LAS VARIABLES

Para el caso que se pretende resolver (un laberinto). Se identifican las siguientes variables para el proceso de aprendizaje automático.

A. Variables de entrada

Que también pueden ser identificadas como variables de observación. Para este caso las entradas del sistema son las siguientes.

- 1) **Posición x**
- 2) **Posición y**

B. Variables de salida

Estas también pueden ser llamadas variables de acción. Dentro de este trabajo se reconocen 2 variables simples.

- 1) **Mover a la derecha**
- 2) **Mover a la izquierda**
- 3) **Mover hacia arriba**
- 4) **Mover hacia abajo**

IX. CAPTURAS DEL FUNCIONAMIENTO DEL SOFTWARE

Primero tomaremos una captura del juego (Figura 3).

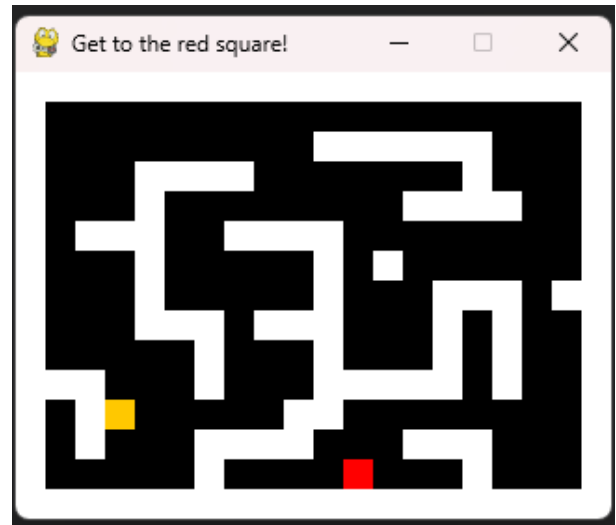


Fig. 3. Muestra del juego en proceso

Segundo mostraremos un link donde podrá visualizar el funcionamiento de programa mediante un video

Ahora los resultados de Q son mostrados en las figuras 4 y 5.

```

✓ 1m 6.6s
Output exceeds the size limit. Open the full output data in a text editor
pygame 2.1.2 (SDL 2.0.18, Python 3.9.15)
Hello from the pygame community. https://www.pygame.org/contribute.html
[[-1.67308090e+01 -1.66958046e+01 -1.67308090e+01 -1.57336456e+01]
 [-1.64805010e+01 -4.11882471e+00 -1.64709755e+01 -1.64411994e+01]
 [-1.61455861e+01 -1.61514900e+01 -1.61867040e+01 -1.29703298e+01]
 [-1.57250533e+01 9.24423768e+00 -1.56903837e+01 -1.56892060e+01]
 [-1.52176291e+01 -7.53896629e-01 -1.52017350e+01 -1.51415693e+01]
 [-1.46217714e+01 -1.45887032e+01 -1.46041033e+01 -1.29932981e+01]
 [-1.41077043e+01 1.69248556e+01 -1.41245833e+01 -1.40744341e+01]
 [-1.35905420e+01 -1.35283293e+01 -1.35581605e+01 -5.89271569e+00]
 [-1.30702658e+01 4.07439960e+01 -1.30875656e+01 -1.30803648e+01]
 [-1.28089542e+01 -1.28089542e+01 -1.82616476e+00 -1.28370864e+01]
 [-1.25468570e+01 -1.25468570e+01 -1.26223053e+01 -1.25754105e+01]
 [-1.21082764e+01 -1.21082764e+01 -1.20904286e+01 -1.17175100e+01]
 [-1.14019680e+01 -1.14019680e+01 -1.14476979e+01 -8.98472630e+00]
 [-1.06899835e+01 -1.06899835e+01 -1.08090701e+01 -1.95139902e+00]
 [-9.97227747e+00 -9.97227747e+00 -1.00360551e+01 1.34752791e+01]
 [-1.14019680e+01 -1.14019680e+01 -1.14476979e+01 -8.98472630e+00]
 [-1.06899835e+01 -1.06899835e+01 -1.08090701e+01 -1.95139902e+00]
 [-9.97227747e+00 -9.97227747e+00 -1.00360551e+01 1.34752791e+01]

```

Fig. 4. Primera parte del resultado de Q

```

[-1.14019680e+01 -1.14019680e+01 -1.14476979e+01 -8.98472630e+00]
 [-1.06899835e+01 -1.06899835e+01 -1.08090701e+01 -1.95139902e+00]
 [-9.97227747e+00 -9.97227747e+00 -1.00360551e+01 1.34752791e+01]
 [-9.24880386e+00 3.46937905e+01 -9.29438799e+00 -9.25449519e+00]
 [-8.79368343e+00 2.88171578e+01 -8.80182415e+00 -8.74127430e+00]
 [-8.70238582e+00 -8.70805132e+00 -2.09343831e+00 -8.70238582e+00]
 [-1.67810690e+01 -1.67789452e+01 -1.68140782e+01 -1.62682292e+01]
 [-1.65164312e+01 -1.65151086e+01 -1.65166595e+01 2.19664434e+01]
 [-1.61520791e+01 -1.62077858e+01 -1.61345717e+01 4.84884164e+01]
 [-1.55861978e+01 -1.55562502e+01 -1.55621165e+01 5.48866535e+01]
 [-1.50079832e+01 -1.50478097e+01 -1.50817958e+01 5.99131816e+01]
 ...
 [-1.99900000e-01 -1.99900000e-01 1.46521315e+01 -1.99900000e-01]
 [ 9.60266581e+00 -3.15089242e+00 -3.15089242e+00 -3.15641785e+00]
 [ 2.19035793e+01 -3.24774153e+00 -3.22110338e+00 -3.16571452e+00]
 [-3.29799868e+00 -3.24774153e+00 -3.29249726e+00 -3.24774153e+00]]

```

Fig. 5. Segunda parte del resultado de Q

X. CONCLUSIONES

- Se concluye que el algoritmo de Q-learning funciona satisfactoriamente según los objetivos de este informe, llegando incluso a obtener el camino más corto.
- Se concluye además que el sistema podría también resolver laberintos mucho más grandes con el número adecuado de episodios.
- Finalmente, se da como una observación que aunque se trata de una posibilidad extremadamente baja, existe la posibilidad de que tal vez el camino encontrado por el sistema sea uno más largo.

XI. CÓDIGO FUENTE

El código fuente empleado se encuentra en el siguiente link <https://drive.google.com/drive/folders/1LF0dUVzh69ejH-xOkXgiasi2vp5oef-O?usp=sharing>

REFERENCES

- [1] J. Torres, *Introducción al aprendizaje por refuerzo profundo*. Independently Published, 2021.
- [2] V. Kanade, "What is reinforcement learning? working, algorithms, and uses." [Online]. Available: <https://www.spiceworks.com/tech/artificial-intelligence/articles/what-is-reinforcement-learning/>
- [3] J. Beakcheol, K. Myeonghwi, and H. Gaspard, *Q-learning Algorithms: A Comprehensive Classification and Applications*, 2017.