

# Aprendizaje por Refuerzo

Pablo Barranco, Thomas Bladt, Alejandro De Anda

Simulación: Proyecto Final

Instituto Tecnológico Autónomo de México

Dr. Jorge de la Vega Góngora

14 de Diciembre del 2020

## Resumen

Tenemos como objetivo mostrar la aplicación de algoritmos de aprendizaje por refuerzo, basados en los métodos de Monte Carlo y Q-Learning. Nos centramos en el problema de la resolución de laberintos, creados mediante una variación del algoritmo de PRIM. Comparamos los métodos, sus soluciones y la eficiencia computacional de cada uno. Finalmente presentamos una breve discusión de las ventajas de cada método.

## 1. Introducción

Este trabajo representa el proyecto final de la materia de Simulación (así como nuestro último trabajo de la licenciatura). Quisimos implementar el uso de inteligencia artificial desde cero (sin utilizar librerías existentes). La mayor parte de los códigos existentes se encuentran programados en `Python`, por lo que decidimos programar nuestros propios códigos en `R`.

Dado que gran parte del material corresponde a temas relacionados con la materia de simulación, pero que no fueron vistos en clase, tuvimos que buscar buenas referencias. Principalmente nos basamos en las clases en línea del Dr. Daniel Soper, “[Foundations of Reinforcement Learning](#).” [3], “[Foundations of Q-Learning](#).” [4], así como en el libro principal de aprendizaje de máquina “*Reinforcement learning: An introduction*” de Richard S. Sutton [1].

## 2. Aprendizaje por Refuerzo

Antes de hablar del concepto de aprendizaje por refuerzo como técnica computacional, podemos entenderlo como un proceso de aprendizaje general. En su contexto más básico, se entiende como un proceso de aprendizaje mediante recompensas y/o castigos. Por ejemplo, pensemos en cómo un niño pequeño aprende a hacer algo nuevo. El niño observa su entorno, toma una decisión y actúa de determinada forma. Posteriormente observa cómo repercutió esa acción en su entorno y aprende de ello. Siguiendo con el ejemplo, si el niño consigue comer solo, sus padres lo felicitarán (recompensa), mientras que si acerca las manos al fuego, será regañado (castigo). Así, el niño aprende a tomar decisiones de una mejor manera. En ocasiones, se puede tomar una decisión que no maximiza la recompensa inmediata, sino que piensa a futuro para maximizar la recompensa total; tal sería el caso de faltar a una reunión de amigos (recompensa inmediata) para poder estudiar y obtener una buena calificación final (recompensa total).

Para poder traducir problemas de la vida real al lenguaje de una computadora, es necesario hacer un planteamiento teórico. Lo primero que se necesita es un agente explorador que tenga un propósito o meta que quiera lograr, que interactúe con un ambiente y obtenga información acerca de sus acciones en cada estado.

En el campo de aprendizaje de máquina existen tres paradigmas principales de aprendizaje: *supervisado*, *no supervisado* y *por refuerzo*[1]. En aprendizaje supervisado se utilizan datos de entrenamiento, que son ejemplos proporcionados por un supervisor externo, en donde se le indica a la computadora cómo reaccionar ante determinadas situaciones. A través de estos ejemplos, el algoritmo busca encontrar una estructura subyacente para poder reaccionar de manera correcta ante futuros eventos (que no sean de prueba). En aprendizaje no supervisado, de igual forma se intenta encontrar una estructura subyacente en un entorno. Sin embargo, no se proporcionan datos de entrenamiento. La diferencia principal que tiene aprendizaje por refuerzo ante estos dos métodos es que no necesariamente busca encontrar una estructura subyacente en el entorno, sino que se basa en un sistema que busca maximizar las recompensas obtenidas.

### 2.1. Entrada-Salida

Como en cualquier algoritmo, aprendizaje por refuerzo requiere de dos cosas: una entrada (input) y una salida (output). Como entrada se le tienen que proporcionar al algoritmo los *estados* del entorno, así como las *recompensas* de estar en cada estado. Como salida, el algoritmo nos dará acciones. Lo que se busca con el algoritmo es determinar una estrategia (policy) óptima para maximizar las recompensas acumuladas. Es decir, busca poder elegir una acción - en cada estado - de manera óptima.

## 2.2. Recompensas

Aprendizaje por refuerzo se basa en un sistema de recompensas. Una recompensa se puede interpretar como una métrica que nos dice qué tan buena resulta ser una acción. Hay que tener en cuenta que lo que busca el algoritmo es maximizar la recompensa acumulada. Es posible que se elija una acción que no maximice la recompensa inmediata, sino una que maximice la recompensa en un futuro. Es por esto que no se le puede considerar a este algoritmo como uno codicioso. Matemáticamente, se define una función recompensa, que dependerá del problema y del modelo (puede querer maximizar ganancias o minimizar pérdidas).

## 2.3. Entorno

El entorno es el “lugar” en donde el agente se encuentra aprendiendo. Es aquello que le provee al algoritmo la información acerca de los estados en los que se encuentra y las recompensas correspondientes. Notamos que no necesariamente tiene que ser interpretado de forma espacial, sino sólo como un ente que provee de información al agente. Por ejemplo, un laberinto puede ser particionado en una cuadrícula, en donde cada cuadro puede ser visto como un estado; en el juego de blackjack un estado dependería de las cartas que se muestran en la mesa y las que ya han salido.

Una vez que se tienen estados, es natural que para cada uno de ellos se tenga una lista finita de posibles acciones. Podría ser una acción espacial (moverse hacia alguna dirección) u otro tipo de acción (en el caso de blackjack podría ser tomar una nueva carta o no). El entorno también define cómo se transiciona de un estado a otro dependiendo de la acción elegida por el agente. El entorno también tiene la información de la recompensa que se tiene por elegir cualquier acción en cualquier estado.

Conforme el algoritmo va recopilando información de las recompensas en cada estado, va aprendiendo. Algo importante que debemos notar es que el agente a veces decide elegir lo que conoce como la mejor acción (que aprendió anteriormente) o una acción completamente aleatoria. Esto permite al agente “explorar” distintas alternativas y, posiblemente, encontrar una mejor solución al problema. La probabilidad de explorar o elegir la mejor acción se controla con un parámetro que define qué proporción del tiempo el agente explora y qué proporción del tiempo elige la mejor acción.

## 2.4. Proceso de Decisión de Markov

Los MDP (procesos de decisión de Markov, por sus siglas en inglés) son el lenguaje matemático mediante el cual el agente decidirá qué acción elegir en el proceso de aprendizaje. Este tipo de sistema estocástico se utiliza cuando una decisión depende en parte de un factor aleatorio y en parte de un ente que toma decisiones (en este caso el agente).

Una parte importante de estos sistemas es que trabajan con tiempos discretos (es decir, cada unidad de tiempo se puede definir como  $t_0, t_1, \dots$ ). El proceso de decisión, al tiempo  $t$ , es como sigue:

1. Entrar al estado  $S_t$ .
2. Elegir una acción  $a_t$ .
3. Recibir una recompensa  $r_t$  por la acción.

con lo que se entra al estado  $S_{t+1}$  y se repiten estos pasos hasta que termine el proceso.

## 2.5. Entrenamiento

Al igual que en aprendizaje supervisado, la primera etapa de aprendizaje por refuerzo consiste en entrenar al algoritmo. Esto se hace mediante muchas iteraciones de transición de estados, tomando en cuenta las recompensas observadas. Se corre el “juego” (entrecomillamos juego, pues lo que se está entrenando no necesariamente es con fines de diversión) repetidas ocasiones, de tal forma que se intente lograr la meta final. Después de cada iteración el algoritmo tiene más información, por lo que actualiza su política de decisión. Una vez que el método haya sido entrenado satisfactoriamente, el sistema está listo para implementarse en un ejemplo real. Es decir, la fase de entrenamiento le sirvió al algoritmo para saber cómo reaccionar ante ciertas situaciones, y ahora es momento de afrontarse por si solo a su tarea. A esta etapa se le conoce como modo de inferencia, y aquí ya no se actualiza la política de decisión, sino que simplemente ejecuta la estrategia que generó en el entrenamiento.

## 3. Creación de los laberintos

Para poder implementar los algoritmos de aprendizaje por refuerzo, lo primero que requerimos fue un ambiente o entorno. Aunque existen bibliotecas de juegos en algunos lenguajes de programación (por ejemplo `pygame` en `Python`), optamos por crear un juego sencillo - la resolución de un laberinto - para tener más control sobre el código.

Crear laberintos a mano es posible, pero quisimos poder generarlos de forma arbitraria y del tamaño de nuestra elección. Para esto, programamos una variante del algoritmo de PRIM; el código está basado en el [pseudocódigo encontrado en Stack Overflow \[2\]](#). Para entenderlo, se necesita saber lo siguiente: el laberinto es una matriz llena de caminos (1) y paredes (0). Una celda frontera es cualquier celda con distancia 2 a la actual, tal que sea pared. Una celda vecina es una celda con distancia 2 a la actual tal que sea camino. A continuación se muestra el pseudocódigo.

---

**Algorithm 1** PRIM modificado

---

**Datos:** tamaño del laberinto  $m, n$ .

**Resultado:** matriz de  $m \times n$  correspondiente a un laberinto.

crear una matriz de  $m \times n$  llena de ceros (paredes)

elegir una celda aleatoria, cambiarla a 1 (camino) y calcular las celdas frontera

**mientras** *lista de celdas frontera no vacía* **hacer**

    elegir celda frontera aleatoria de la lista y hacerla celda actual

    elegir celda vecina aleatoriamente y conectarla con la celda actual

    calcular celdas frontera de la actual y agregarlas a la lista

    quitar celda actual de la lista

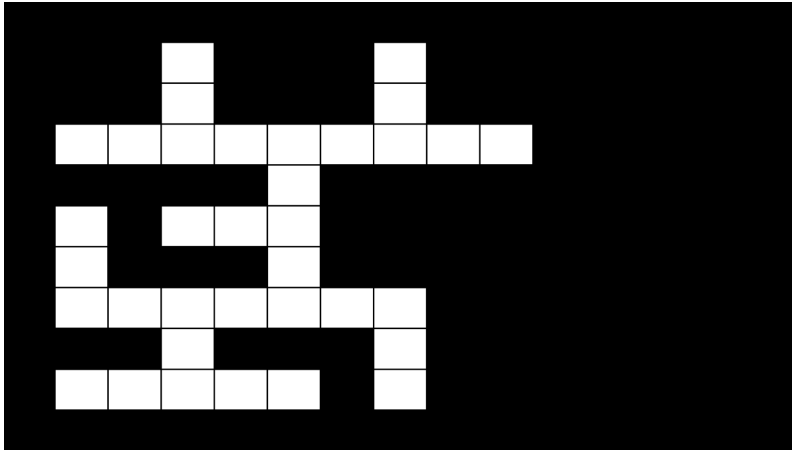
**fin**

---

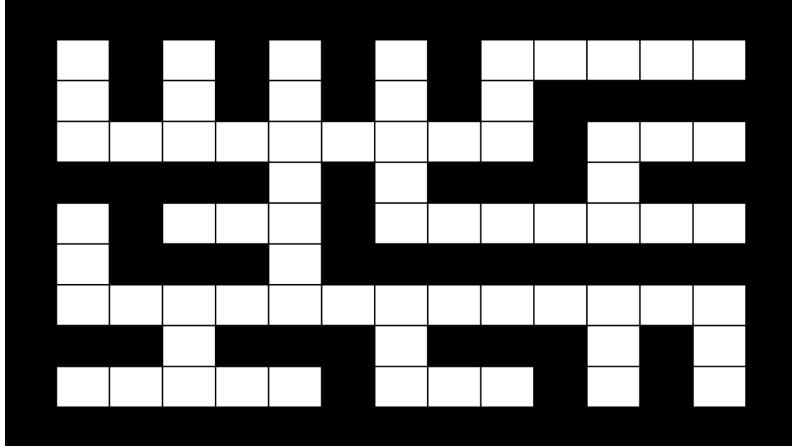
En la Figura [1](#) mostramos distintas etapas de la creación de un laberinto de tamaño  $11 \times 15$ , utilizando nuestro algoritmo.



(a) Primera iteración.



(b) Iteración número 17.



(c) Iteración final.

Figura 1: Distintas etapas del algoritmo de PRIM modificado.

Finalmente, elegimos como entrada a la celda  $(m - 1, 1)$  y salida a  $(2, n)$ . Teniendo ya un laberinto, se crea la matriz de recompensas. Una vez desarrollado el laberinto, debemos construir su matriz de recompensas. Esta matriz consta de valores negativos únicamente, excepto en la celda correspondiente a la meta (cuya recompensa es positiva y muy grande).

Además, los valores son más grandes en las celdas en las que existen paredes, mientras que el valor es más pequeño en las celdas que contienen un camino. Esto es porque queremos disuadir al algoritmo de quedarse en un bucle infinito en el laberinto, maximizando cada vez más la recompensa; al asignarle un castigo a cada unidad de camino, el algoritmo buscará llegar a la recompensa de la forma más eficiente posible. En el cuadro 1 podemos ver la matriz de recompensas resultante de nuestro ejemplo y en la figura 2 mostramos la matriz de una forma más visual.

-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100
-100	-1	-100	-1	-100	-1	-100	-1	-100	-1	-1	-1	-1	-1	1000
-100	-1	-100	-1	-100	-1	-100	-1	-100	-1	-100	-100	-100	-100	-100
-100	-1	-1	-1	-1	-1	-1	-1	-1	-1	-100	-1	-1	-1	-100
-100	-100	-100	-100	-100	-1	-100	-1	-100	-100	-100	-1	-100	-100	-100
-100	-1	-100	-1	-1	-1	-100	-1	-1	-1	-1	-1	-1	-1	-100
-100	-1	-100	-100	-100	-1	-100	-100	-100	-100	-100	-100	-100	-100	-100
-100	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-100
-100	-100	-100	-1	-100	-100	-100	-1	-100	-100	-100	-1	-100	-1	-100
-100	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-100
-100	-100	-100	-1	-100	-100	-100	-1	-100	-100	-100	-1	-100	-1	-100
-1	-1	-1	-1	-1	-1	-100	-1	-1	-1	-100	-1	-100	-1	-100
-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100	-100

Cuadro 1: Matriz de recompensas para el laberinto.

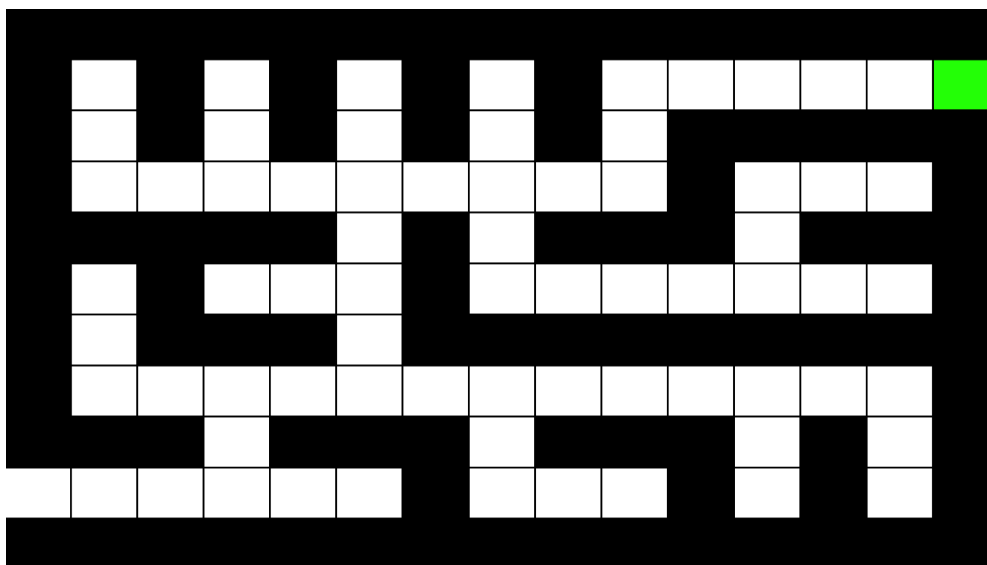


Figura 2: Laberinto final. La celda verde corresponde a la meta.

## 4. Q-Learning

Q-Learning es un tipo particular de aprendizaje por refuerzo, por lo que contamos con un agente explorador y con un entorno en el que éste se encuentra. A diferencia de otros métodos, en Q-Learning se busca construir una conducta mediante la interacción directa con el entorno. Así, el agente explorador interactúa una y otra vez con el entorno

y aprende mediante un mecanismo de prueba/error. Al final de cada intento, el agente actualiza sus conocimientos y es así como aprende a desarrollarse en su entorno.

Q-Learning cuenta con todas las características de un modelo de aprendizaje por refuerzo. Sin embargo, Q-Learning también requiere que tanto el conjunto de estados como de acciones sean finitas. Es decir, existe un conjunto finito de estados  $\mathcal{S} = \{s_0, s_1, \dots, s_n\}$  en los que el agente puede encontrarse y en cada uno de ellos puede elegir entre un conjunto finito de acciones  $\mathcal{A} = \{a_0, a_1, \dots, a_m\}$ . Considerando nuestro ejemplo, supongamos que tenemos un laberinto de tamaño  $10 \times 10$ ; entonces, cada celda corresponde a un estado, con lo que  $\mathcal{S} = \{1, 2, \dots, 100\}$  y las acciones que puede tomar el agente son moverse arriba, derecha, izquierda o abajo, con lo que podemos definir  $\mathcal{A} = \{1, 2, 3, 4\}$ .

Este método debe su nombre a que funciona mediante *Q-values*. Estos valores representan la calidad (quality en inglés, de ahí la Q) de tomar la decisión  $a$  en el estado  $s$ , que podemos escribir como una función  $Q(s, a)$ . Intuitivamente, un Q-value nos dice qué tanto esperamos obtener de recompensa acumulada al tomar una acción en determinado estado.

Regresando a nuestro ejemplo, necesitamos una forma de guardar los Q-values. Para cada estado, necesitamos un Q-value para cada acción. Por lo tanto, requerimos de un array de dimensión  $m \times n \times 4$ , donde  $m$  representa el número de filas del laberinto y  $n$  el número de columnas. La idea es que cuando el agente se encuentre en algún estado, elija la acción que tenga un mayor Q-value. De esta forma, se puede interpretar al array de los Q-values como la política del agente.

## 4.1. Diferencias Temporales

Ahora debemos pensar en cómo se actualizarán los Q-values conforme el agente va recopilando información de su entorno. En concreto: ¿Cómo actualizar los Q-values de determinada acción, en determinado estado, después de haberla realizado? Para esto existe un método conocido como diferencias temporales, que nos proporciona la siguiente ecuación

$$DT(s_t, a_t) = r_t + \gamma \cdot \max_a Q(s_{t+1}, a) - Q(s_t, a_t). \quad (1)$$

En la ecuación 1 tenemos que

- $r_t$  es la recompensa inmediata (dada por la matriz de recompensas) por haber tomado la acción  $a_t$ .
- $\gamma$  es un parámetro de descuento, entre 0 y 1, que controla qué tan importantes son las recompensas en un futuro. Es intuitivo que una recompensa inmediata es preferible que una futura, y  $\gamma$  nos dice qué tanto.



- $\max_a Q(s_{t+1}, a)$  corresponde al Q-value más grande en el estado actual  $s_{t+1}$ . Nos dice cuánto esperamos acumular en recompensas futuras, en el mejor de los casos.
- $Q(s_t, a_t)$  es el Q-value en el estado pasado, con la acción elegida. Es el valor que buscamos actualizar.

En resumen, la ecuación 1 nos permite considerar recompensas inmediatas y futuras a la vez.

## 4.2. Ecuación de Bellman

Una vez que tenemos la diferencia temporal calculada, es hora de actualizar el Q-value. Para esto introducimos la ecuación de Bellman, dada por

$$Q_{\text{nueva}}(s_t, a_t) = Q_{\text{vieja}}(s_t, a_t) + \alpha \cdot DT(s_t, a_t). \quad (2)$$

En la ecuación 2 tenemos que

- $Q_{\text{nueva}}(s_t, a_t)$  corresponde al valor actualizado del Q-value en el estado  $s_t$  realizando la acción  $a_t$
- $Q_{\text{vieja}}(s_t, a_t)$  corresponde al valor pasado del Q-value en el mismo estado y bajo la misma acción.
- $\alpha$  es un parámetro de aprendizaje, entre 0 y 1, que controla qué tan rápido se actualizan los Q-values, basado en la diferencia temporal.
- $DT(s_t, a_t)$  es la diferencia temporal, calculada en la ecuación 1.

Esta ecuación es sumamente importante dado que actualiza de manera adecuada los Q-values, dependiendo de la acción que se tomó en el estado anterior. Como su nombre lo indica, esta ecuación fue inventada por Richard Bellman.

Entonces, recapitulando: el array de Q-values se inicializa como un array de ceros (puede ser lo que sea, en realidad), y en cada iteración se realiza el proceso de pasar de un estado a otro, a través de las acciones posibles. En cada paso se calcula la diferencia temporal, la ecuación de Bellman y se actualiza el Q-value. En la fase de entrenamiento se repite el proceso muchas veces, con lo que el array de Q-values se va refinando cada vez más. Una vez que concluye la fase de entrenamiento, obtenemos un array de Q-values que refleja la política que se debe seguir para maximizar la recompensa. Qué tan precisa es esa política depende de varios factores, uno de los más grandes siendo el número de iteraciones para el entrenamiento (problemas más complejos requieren más entrenamiento).

### 4.3. Implementación del método

Para programar el método de Q-Learning nos basamos [en un código de Python para un problema similar](#) [5], por parte del Dr. Daniel Soper. Presentamos el pseudocódigo de nuestra implementación:

---

**Algorithm 2** Q-Learning.

---

**Datos:** número de entrenamientos, matriz de recompensas,  $\varepsilon$ ,  $\alpha$  y  $\gamma$ .

**Resultado:** array de Q-values.

```
para  $i = 1$  hasta número de entrenamientos hacer
    elegir celda de inicio aleatoria.
    mientras no sea una celda terminal hacer
        elegir siguiente acción
        calcular siguiente estado (nueva celda)
        calcular recompensa inmediata
        calcular diferencia temporal
        actualizar Q-value
    fin
fin
```

---

## 5. Monte Carlo

Este algoritmo es más primitivo que Q-learning. La diferencia principal radica en que no utiliza la ecuación de Bellman para mapear las recompensas, sino que es un simple promedio. El término “Monte Carlo” por lo general se emplea para modelos de estimación que están compuestos por un término aleatorio. En este caso, sí existe el uso de un componente aleatorio para la exploración del agente llamado “ $\varepsilon$  - greedy”. Pero no por eso este algoritmo lleva el nombre de Monte Carlo, sino porque realiza un promedio de los retornos del entorno.

Monte Carlo, al igual que Q-Learning, requiere que tanto el conjunto de estados como el conjunto de acciones sean finitos. Tras múltiples simulaciones, el agente interactúa con el entorno de manera finita y estas interacciones reciben el nombre de “episodios”. De forma iterativa, los episodios consiguen que el agente aprenda a interactuar de manera óptima dentro del entorno y de esta forma maximizan la recompensa previamente presentada.

Para entender Monte Carlo, primero hay que entender la forma de exploración del algoritmo. De manera similar a Q-Learning, este método explora el entorno si un número aleatorio es menor a  $\varepsilon$ . De lo contrario, el agente selecciona la siguiente acción que maximice el retorno esperado. De esta forma, a mayor  $\varepsilon$ , mayor exploración. En pocas pa-

labras, lo que hace el algoritmo de Monte Carlo es seleccionar un estado inicial aleatorio. Mediante el uso de  $\varepsilon$  - greedy, selecciona una nueva acción y, consecuentemente, un nuevo estado. Finalmente, repite esto hasta llegar a un estado terminal. Una vez finalizado el episodio, se promedia la recompensa obtenida tras la primer visita de cada estado. Así, se registra memoria de los estados visitados y sus recompensas asociadas multiplicadas por un factor de descuento. Esto favorece a que el aprendizaje del agente sea robusto, priorizando los retornos inmediatos sin olvidar las últimas acciones del episodio.

La asignación del retorno ponderado al estado del episodio en el tiempo  $t$  se asigna mediante la siguiente ecuación:

$$G_t = R_{t+1} + \gamma R_{t+1} + \gamma^2 R_{t+3} + \dots,$$

donde  $0 \leq \gamma \leq 1$  es la tasa de retorno y  $R_t$  es el retorno en el tiempo  $t$ .

De esta forma, se aproxima el retorno esperado en cada estado. Para posteriormente optar por la política cuyas acciones maximicen el retorno.

A continuación se presentan los pseudocódigos de nuestra implementación de Monte Carlo:

---

**Algorithm 3** Monte Carlo: estructura general.

---

**Datos:** número de entrenamientos, matriz de recompensas,  $\varepsilon$  y  $\gamma$ .

**Resultado:** array de Q-values.

---

Inicializar  $Q$ ,  $v$  = visitas y  $r$  = recompensas como arrays de ceros.

**para**  $i = 1$  *hasta número de entrenamientos* **hacer**

$\varepsilon \leftarrow$  máximo entre 0 y  $\varepsilon - i/\text{número de entrenamientos}$

episodio  $\leftarrow$  episodio( $\varepsilon, \gamma$ )

inicializar vector de estado-acción encontrados

**para**  $k = 1$  *hasta duración de episodio* **hacer**

**si** estado-acción de episodio no se ha encontrado **entonces**

$v(\text{estado}, \text{acción}) \leftarrow v(\text{estado}, \text{acción}) + 1$

$r(\text{estado}, \text{acción}) \leftarrow r(\text{estado}, \text{acción}) + r$

$Q(\text{estado}, \text{acción}) \leftarrow r(\text{estado}, \text{acción})/v(\text{estado}, \text{acción})$

cambiar estado-acción a encontrado

**fin**

**fin**

**fin**

---

---

**Algorithm 4** Monte Carlo: episodio.

---

**Datos:**  $\varepsilon$  y  $\gamma$ .**Resultado:** tabla estado-acción-recompensa.

inicializar tabla estado-acción-recompensa como tabla vacía de 3 columnas

elegir estado inicial aleatorio

**mientras** *estado no terminal* **hacer**| elegir siguiente acción  $a$ 

| calcular siguiente estado (nueva celda)

| calcular recompensa inmediata  $r$ 

| agregar a estado-acción-recompensa la fila (estado anterior, acción, recompensa)

**fin** $g \leftarrow 0$  $R \leftarrow$  estado-acción-recompensa al revés**para**  $i = 1$  *hasta* *dimensión*  $R$  **hacer**|  $g \leftarrow r + \gamma g$ | guardar valores de  $g$  en vector  $G$ **fin**reemplazar vector  $G$  en última columna de estado-acción-recompensa

---

## 6. Resultados

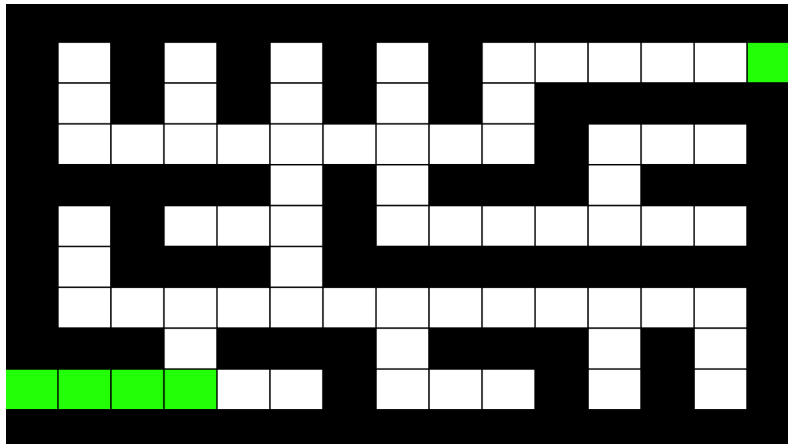
En esta sección presentamos los resultados que obtuvimos para el problema de la resolución de laberintos de distintas dimensiones. Corremos el algoritmo (sea Q-learning original o aquel basado en Monte Carlo) para obtener una matriz  $Q$  y procedemos a la etapa de inferencia, es decir, vemos si el algoritmo puede resolver el laberinto.

### 6.1. Ejemplo práctico

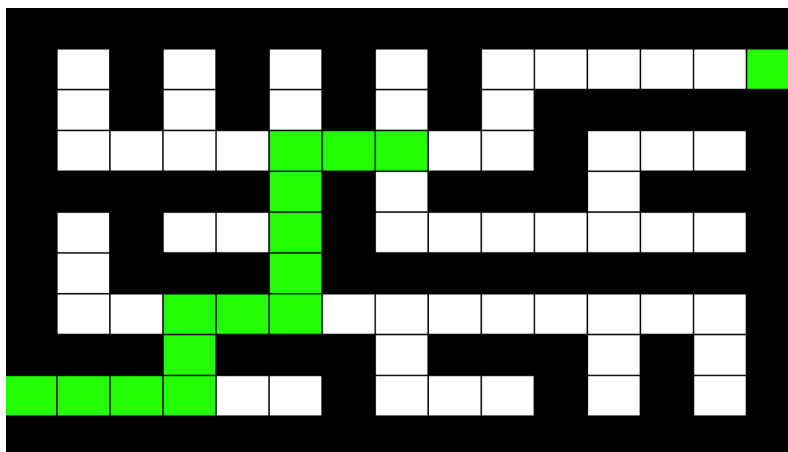
Aquí mostramos todo el proceso para resolver el problema de un laberinto en particular. Lo primero que debemos hacer es crear un laberinto, que elegimos que sea de tamaño  $11 \times 15$ . En particular, vamos a trabajar con el laberinto creado en la sección 3. Corrimos el algoritmo de Q-Learning, entrenando un total de 50,000 iteraciones (se eligió este número a base de prueba y error). El algoritmo tardó un total de 2.832 segundos. El camino que resulta es el siguiente

$[10, 1] \rightarrow [10, 2] \rightarrow [10, 3] \rightarrow [10, 4] \rightarrow [9, 4] \rightarrow [8, 4] \rightarrow [8, 5] \rightarrow$   
 $[8, 6] \rightarrow [7, 6] \rightarrow [6, 6] \rightarrow [5, 6] \rightarrow [4, 6] \rightarrow [4, 7] \rightarrow [4, 8] \rightarrow [4, 9] \rightarrow$   
 $[4, 10] \rightarrow [3, 10] \rightarrow [2, 10] \rightarrow [2, 11] \rightarrow [2, 12] \rightarrow [2, 13] \rightarrow [2, 14] \rightarrow [2, 15]$

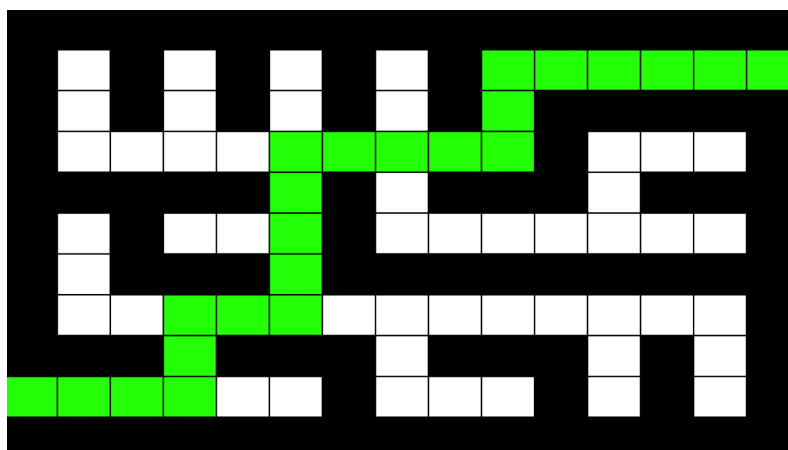
De igual forma lo podemos visualizar gráficamente. En la figura 3 se muestra el camino solución en distintas etapas.



(a) Primera iteración.



(b) Iteración número 17.



(c) Iteración final.

Figura 3: Distintas etapas del algoritmo de PRIM modificado.

Ahora que hemos visto cómo funciona el algoritmo, pasaremos a evaluar la calidad de los métodos y los compararemos.

## 6.2. Evaluación de los métodos

Para evaluar los métodos, corrimos el algoritmo para laberintos cuadrados, de tamaños:

- Laberinto 1:  $5 \times 5$  (25 casillas)
- Laberinto 2:  $7 \times 7$  (49 casillas)
- Laberinto 3:  $9 \times 9$  (81 casillas)
- Laberinto 4:  $11 \times 11$  (121 casillas)
- Laberinto 5:  $13 \times 13$  (169 casillas)
- Laberinto 6:  $15 \times 15$  (225 casillas)

Notamos que los tamaños siempre consisten de números impares. Esto es porque el algoritmo de PRIM modificado genera laberintos simétricos (márgenes simétricos) para estos casos.

Queremos aclarar que el número de iteraciones lo elegimos de manera “artesanal”, basado en nuestras experiencias de prueba y error. Intentamos elegir el menor número posible de iteraciones para el entrenamiento, de tal forma que nos proveyera con una matriz de Q-values con una solución correcta al problema.

## 6.3. Comparativa

Los valores obtenidos con ambos métodos se reportan a continuación en los cuadros 2 y 3.

Cuadro 2: Resultados utilizando Q-Learning.

Laberinto	Tamaño	Casillas	Iteraciones	Tiempo
1	$5 \times 5$	25	1,000	0.151
2	$7 \times 7$	49	2,500	0.240
3	$9 \times 9$	81	12,500	0.647
4	$11 \times 11$	121	12,500	0.820
5	$13 \times 13$	169	25,000	1.232
6	$15 \times 15$	225	50,000	2.819

Cuadro 3: Resultados utilizando Monte Carlo.

Laberinto	Tamaño	Casillas	Iteraciones	Tiempo
1	$5 \times 5$	25	1,000	0.54
2	$7 \times 7$	49	10,000	8.717
3	$9 \times 9$	81	-	-
4	$11 \times 11$	121	-	-
5	$13 \times 13$	169	-	-
6	$15 \times 15$	225	-	-

**Nota.** El método de Monte Carlo no encontró la solución en un tiempo razonable de entrenamiento para los laberintos 3, 4, 5 y 6. Hablamos más al respecto en la sección 7.

Finalmente, quisimos investigar la complejidad del algoritmo. Es claro que un laberinto más grande requiere más iteraciones de entrenamiento, por lo que investigamos cómo crece el tiempo con respecto a las iteraciones. Corrimos el código para laberintos de tamaño  $n = \{50000, 75000, \dots, 500000\}$  para el método de Q-Learning. En la figura ?? se muestran los resultados obtenidos.

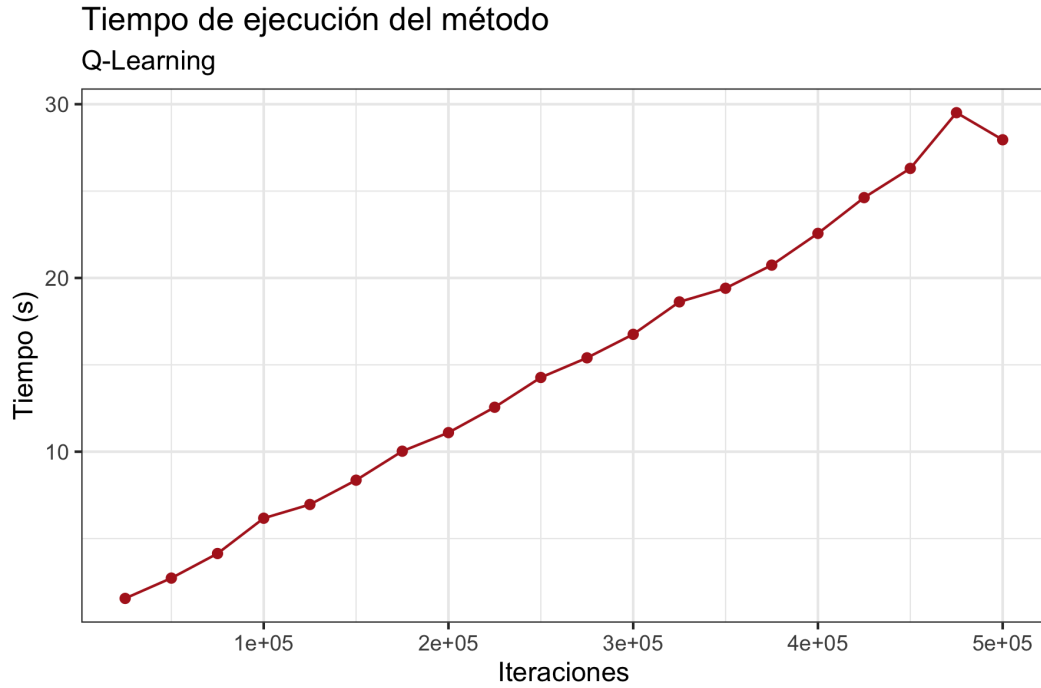


Figura 4: Tiempos de ejecución del método Q-Learning.

Como podemos ver, parecer ser que el algoritmo crece de forma lineal, es decir, concluimos empíricamente que el método tiene orden  $\mathcal{O}(n)$ .

## 7. Conclusiones

Logramos implementar el método de aprendizaje por refuerzo de Q-learning en R, así como el algoritmo de PRIM modificado para la creación de laberintos.

El método de Q-Learning presenta una manera aceptable de tratar el problema de la resolución de los laberintos. En tan solo 2.8 segundos, y después de 50,000 iteraciones, logra aprender satisfactoriamente el ambiente de un laberinto de  $15 \times 15$  y llega a la solución correcta. Pudimos observar que el orden de complejidad de este algoritmo es lineal.

En cuanto a la implementación de un método de Monte Carlo, no resultó ser muy exitosa. Es muy probable que exista un error en el código, pues aún cuando la teoría nos dice que el método es lento (y peor que Q-Learning), el desempeño del método es demasiado malo. A partir del laberinto de tamaño  $9 \times 9$ , el algoritmo no logró llegar a la solución, incluso después de 5,000,000 de iteraciones (el laberinto pasado requirió de tan sólo 10,000). Además, el tensor  $Q$  no parece ser correcto, al inspeccionarlo cuidadosamente. Sin embargo, presentamos la teoría de Monte Carlo, junto con los pseudocódigos, esperando poder corregir el código en un futuro y así lograr una comparación adecuada entre los métodos.



## Referencias

- [1] Sutton, R. S., and Barto, A. G.: *Reinforcement learning: An introduction*. The MIT Press, Cambridge, MA, 2018.
- [2] BitTickler (<https://stackoverflow.com/users/2225104/bittickler>), Implementing a randomly generated maze using Prim's Algorithm, URL (consultado: 14-12-2020): <https://stackoverflow.com/questions/29739751/implementing-a-randomly-generated-maze-using-prims-algorithm>
- [3] Soper, Daniel. "Foundations of Reinforcement Learning.", YouTube, publicado 7 de abril del 2020, URL (consultado: 14-12-2020) ) [www.youtube.com/watch?v=wVXXLLT6srY](http://www.youtube.com/watch?v=wVXXLLT6srY).
- [4] Soper, Daniel. "Foundations of Q-Learning.", YouTube, publicado 22 de abril del 2020, URL (consultado: 14-12-2020) [https://www.youtube.com/watch?v=\\_\\_t2XRxXGxI](https://www.youtube.com/watch?v=__t2XRxXGxI).
- [5] Soper, Daniel. "Notebook for Topic 08 Video - Q-Learning - A Complete Example in Python.ipynb." Google Drive, URL (consultado: 14-12-2020) [colab.research.google.com/drive/1E2RViy7xmor0mhqskZV14\\_NUj2jMpJz3](https://colab.research.google.com/drive/1E2RViy7xmor0mhqskZV14_NUj2jMpJz3).

## Apéndice: código principal

```
1 ---
2 title: "Aprendizaje por Refuerzo"
3 author: "Pablo Barranco, Thomas Bladt, Alejandro De Anda"
4 date: "14/12/2020"
5 output: html_document
6 ---
7
8 ```{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10 rm(list=ls())
11 source("prim_final.R")
12 ```
13
14 # Código para encontrar la ruta más corta entre dos puntos
15 # Aplicado para encontrar la solución a un laberinto
16 # Basado en explicación y código en Python de Dr. Daniel Soper:
17 # https://www.youtube.com/watch?v=iKdlKYG78j4
18 # 14 de diciembre del 2020
19 # Pablo Barranco, Thomas Bladt, Alejandro De Anda
20
21 Primero generamos un laberinto con el algoritmo de PRIM y creamos el
22 ambiente del "juego".
23
24 Definimos función para crear matriz de recompensas y graficarla.
25 ```{r}
26 creaRecompensas <- function(lab, rmax){
27   recompensas <- lab
28   recompensas[recompensas == 0] <- -100
29   recompensas[recompensas == 1] <- -1
30   recompensas[2, ncol(lab)] <- rmax
31   recompensas[nrow(lab)-1, 1] <- -1
32   return(recompensas)
33 }
34
35 plotRecompensas <- function(rec){
36   rec[2, ncol(rec)] <- 100
37   plot(rec, col = c("black", "white", "red", "green"), key = NULL, breaks =
38     NULL, axis.col=NULL, axis.row=NULL, xlab = "", ylab="", main = paste0
39     ("Laberinto de ", nrow(rec), " x ", ncol(rec)))
40 }
41 ```
42
43 Generamos un laberinto de tamaño $m \times n$.
44 ```{r}
45 # Tamaño del laberinto (m x n)
```

```

43 m <- 9
44 n <- 9
45
46 set.seed(2022)
47 laberinto <- genera_laberinto(m,n,0)
48
49 recompensas <- creaRecompensas(laberinto,10000)
50
51 # Inicializamos matriz Q
52 q_val <- array(0, dim = c(m,n,4))
53
54 # Ploteo laberinto
55 #png("lab_15.png",width = 1800, height = 1300,res=250)
56 plotRecompensas(recompensas)
57 #dev.off()
58 '''
59
60 Definimos las funciones ocupadas
61
62 '''{r}
63 es_terminal <- function(fila,columna){
64   ifelse(recompensas[fila,columna] == -1,FALSE,TRUE)}
65
66 inicio_aleatorio <- function(){
67   fila <- sample(1:m,size = 1)
68   columna <- sample(1:n,size = 1)
69   while(es_terminal(fila,columna)){
70     fila <- sample(1:m,size = 1)
71     columna <- sample(1:n,size = 1)
72   }
73   return(c(fila,columna))
74 }
75
76 sig_accion <- function(fila,columna,epsilon){
77   if(runif(1) < epsilon )
78     return(which.max(q_val[fila,columna,]))
79   else
80     return(sample(1:4,size = 1))
81 }
82
83 sig_lugar <- function(fila,columna,accion){
84   if(accion == 1 && fila > 1){fila = fila - 1}
85   else
86     if(accion == 2 && columna < n){columna = columna + 1}
87     else
88       if(accion == 3 && fila < m){fila = fila + 1}
89       else

```

```

90     if(accion == 4 && columna > 1){columna = columna - 1}
91     return(c(fila,columna))
92 }
93
94 dist_mas_corta <- function(fila_inicio,columna_inicio){
95     if(es_terminal(fila_inicio,columna_inicio)){return("[,]")}
96     else
97         fila <- fila_inicio
98         columna <- columna_inicio
99         camino <- paste0("[",fila,",",columna,"]")
100         k <- 0
101
102         solucion_df <- data.frame(x= fila_inicio,y = columna_inicio)
103
104         while(!(es_terminal(fila,columna)) && (k < 500)){
105             k <- k+1
106             accion <- sig_accion(fila, columna,1)
107             aux <- sig_lugar(fila,columna,accion)
108             fila <- aux[1]
109             columna <- aux[2]
110             camino <- append(camino,paste0("[",fila,",",columna,""])
111             solucion_df[k+1,1] <- fila
112             solucion_df[k+1,2] <- columna
113         }
114         return(list(camino = camino,sol = solucion_df))
115     }
116
117 look <- function(vec,mat){
118     resp <- apply(mat, 1, function(x, want) isTRUE(all.equal(x, want)),
119         vec)
120     if(TRUE %in% resp)
121         return(TRUE)
122     else
123         return(FALSE)
124 }
125
126 El algoritmo de Monte Carlo:
127 '''{r}
128 episodio <- function(epsilon,gamma = 0.9) # corre un episodio de
129     iteracion pero con tasa de descuento
130 {
131     bandera = TRUE
132     # Obtenemos S_0
133     inicio <- inicio_aleatorio()
134     fila <- inicio[1]
135     columna <- inicio[2]

```

```

135
136 state_action_reward = c(NULL, NULL, NULL)
137
138 while(!(es_terminal(fila, columna))){
139   accion <- sig_accion(fila, columna, epsilon)
140   fila_vieja <- fila
141   columna_vieja <- columna
142   lugar <- sig_lugar(fila, columna, accion)
143   fila <- lugar[1]
144   columna <- lugar[2]
145   recompensa <- recompensas[fila, columna]
146   state_action_reward <- rbind(state_action_reward, c(fila_vieja,
147     columna_vieja, accion, recompensa))
148 }
149
150 g = 0
151 G <- NULL
152 R <- apply(t(state_action_reward[,4]), 1, rev)
153 for(r in R){
154   g = r + gamma*g
155   G <- append(G, g)
156 }
157 G <- apply(t(G), 1, rev)
158 state_action_reward[,4] <- G
159 return(state_action_reward)
160 }
161
162 monte_carlo <- function(env, epsilon, gamma, N_episodes=1000)
163 {
164   # Inicializamos Q, V, PI
165   dim <- dim(env)
166   m <- dim[1]
167   n <- dim[2]
168   q_val <- array(0, dim = c(m, n, 4))
169   visit <- array(0, dim = c(m, n, 4))
170   rewards <- array(0, dim = c(m, n, 4))
171
172   for(i in 1:N_episodes){
173
174     #cambiar epsilon
175
176     eps = max(0, epsilon - i/N_episodes)
177     #generamos un episodio
178     episode <- episodio(eps, gamma) #aqui
179     seen_state_action = matrix(nrow=1, ncol=3)
180     k <- dim(episode)[1]

```

```

181
182     for(i in 1:k){
183         if(!look(episode[i,1:3],seen_state_action)){ #Si no esta en los
184             estados vistos
185             # Estados
186             fila <- episode[i,1]
187             columna <- episode[i,2]
188             # Accion
189             accion <- episode[i,3]
190             #Recompensa
191             recompensa <- episode[i,4]
192
193             #Actualizacion
194             visit[fila, columna, accion] <- visit[fila, columna, accion] +
195             1
196             rewards[fila,columna,accion] <- rewards[fila,columna,accion] +
197             recompensa
198             q_val[fila,columna,accion] <- rewards[fila,columna,accion]/
199             visit[fila, columna, accion]
200             seen_state_action <- rbind(seen_state_action,episode[i,1:3])
201         }
202     }
203     return(q_val)
204 }
205
206 '''
207 '''{r}
208
209 tiempos_x <- seq(from = 25000, to = 500000, by = 25000)
210 tiempos_y <- rep(0,20)
211
212 for(i in 1:20){
213     q_val <- array(0, dim = c(m,n,4))
214     tiempos_y[i] <- system.time(q_val <- entrena_qlearn(tiempos_x[i]))[3]
215 }
216
217 '''
218
219 El algoritmo de Q-Learning:
220 '''{r}
221 entrena_qlearn <- function(num_entren,q = q_val){
222
223     epsilon <- 0.9
224     descuento <- 0.9 #gamma
225     tasa_aprendizaje <- 0.9 #alpha
226
227     for(i in 1:num_entren){

```

```

224 inicio <- inicio_aleatorio()
225 fila <- inicio[1]
226 columna <- inicio[2]
227
228 while(!(es_terminal(fila,columna))){
229     accion <- sig_accion(fila,columna,epsilon)
230     fila_vieja <- fila
231     columna_vieja <- columna
232     lugar <- sig_lugar(fila,columna,accion)
233     fila <- lugar[1]
234     columna <- lugar[2]
235     recompensa <- recompensas[fila,columna]
236     viejo_q_val <- q_val[fila_vieja,columna_vieja,accion]
237     dif <- recompensa + descuento*max(q_val[fila,columna,])-viejo_q_
val
238
239     nuevo_q_val = viejo_q_val + (tasa_aprendizaje*dif)
240     q_val[fila_vieja, columna_vieja, accion] = nuevo_q_val
241 }
242 }
243
244 return(q_val)
245 }
246 '''
247
248 '''{r}
249 #entrena
250 q_val <- array(0, dim = c(m,n,4))
251 system.time(q_val <- entrena_qlearn(100000))
252
253 #Resuelve para la ruta
254 (solucion <- dist_mas_corta(m-1,1))
255
256 #coloreaa
257 colorear <- function(rec,sol){
258     aux <- list()
259     for(i in 1:nrow(sol)){
260         aux[[i]] <- rec
261         rec[sol[i,1],sol[i,2]] <- 100
262     }
263
264     saveGIF({
265         for (i in 1:nrow(sol)) plot( aux[[i]],col = c("black","white","
green"),key = NULL,breaks = NULL,axis.col=NULL, axis.row=NULL,xlab =
"", ylab="", main = "Soluci n del laberinto")
266     },interval = 0.25,movie.name = "solucion1.gif")
267 }

```

```

268
269 #guardas gif coloreado
270 colorea(recompensas,solucion$sol)
271 recompensas[2,49] <- 100
272 '''
273
274 '''{r}
275 #gr fica de complejidad
276 tiempos <- data.frame(x= tiempos_x, y = tiempos_y)
277
278 ggsave("tiempo_q.png",width = 6, height = 4)
279 ggplot(data = tiempos, aes(x= x, y = y))+
280   geom_line(color = "firebrick")+
281   geom_point(color = "firebrick")+
282   labs( x= "Iteraciones", y = "Tiempo (s)", title= "Tiempo de
    ejecuci n del m todo", subtitle = "Q-Learning")+
283   theme_bw()
284 dev.off()
285 '''

```



## Apéndice: código PRIM

```
1 #
   #####
2 ## C digo para generar laberintos utilizando el algoritmo de PRIM
   aleatorio
3 ## Basado en el pseudoc digo encontrado en:
4 ## https://stackoverflow.com/questions/29739751/implementing-a-randomly-generated-maze-using-prims-algorithm
5 ## 19 de noviembre del 2020
6 ## Pablo Barranco, Thomas Bladt, Alejandro De Anda
7 #
   #####
8
9 # Librerias necesarias
   -----
10 library(animation)
11 library(plot.matrix)
12
13 # Funciones auxiliares
   -----
14
15 # Crea la base del laberinto de mxn
16 creaGrid <- function(num_fil,num_col){
17   grid_lab <- matrix(0, nrow = num_fil,ncol = num_col)
18   return(grid_lab)
19 }
20
21 # Selecciona una celda aleatoriamente
22 celdaAleatoria <- function(num_fil,num_col){
23   fila <- sample(2:num_fil-1,size = 1, replace = F)
24   columna <- sample(2:num_col-1,size = 1, replace = F)
25   return(c(fila,columna))
26 }
27
28 # Checa si determinada celda cae dentro del laberinto
29 esLegal <- function(fila,columna, lab){
30   num_fil <- nrow(lab)
31   num_col <- ncol(lab)
32   return( (fila > 1) &&(fila < num_fil) && (columna > 1)&& (columna <
     num_col) )
33 }
34
35 # Calcula la frontera de una celda espec fica
36 # Frontera: celdas a distancia 2 de tipo "cerradas"
```

```

37 frontera <- function(fila,columna,lab){
38   #Inicializamos un data frame
39   frontera_df <- data.frame(x = integer(), y = integer())
40
41   #Consideramos los 4 puntos posibles
42   if(esLegal(fila+2,columna,lab)){
43     if(lab[fila+2,columna] == 0)
44       frontera_df <- rbind(frontera_df,data.frame(x= fila+2, y =
         columna))
45   }
46   if(esLegal(fila-2,columna,lab)){
47     if(lab[fila-2,columna] == 0)
48       frontera_df <- rbind(frontera_df,data.frame(x= fila-2, y =
         columna))
49   }
50   if(esLegal(fila,columna+2,lab)){
51     if(lab[fila,columna+2] == 0)
52       frontera_df <-rbind(frontera_df,data.frame(x= fila, y = columna
         +2))
53   }
54   if(esLegal(fila,columna-2,lab)){
55     if(lab[fila,columna-2] == 0)
56       frontera_df <-rbind(frontera_df,data.frame(x= fila , y = columna
         -2))
57   }
58
59   return(frontera_df)
60 }
61
62 # Calcula las celdas vecinas a una espec fica
63 # Vecina: celdas a distancia 2 de tipo "abierta"
64 vecinos <- function(fila,columna,lab){
65   vecinos_df <- data.frame(x = integer(), y = integer())
66
67   if(esLegal(fila+2,columna,lab)){
68     if(lab[fila+2,columna] == 1)
69       vecinos_df <- rbind(vecinos_df,data.frame(x= fila +2, y = columna
         ))
70   }
71   if(esLegal(fila-2,columna,lab)){
72     if(lab[fila-2,columna] == 1)
73       vecinos_df <-rbind(vecinos_df,data.frame(x= fila-2, y = columna))
74   }
75   if(esLegal(fila,columna+2,lab)){
76     if(lab[fila,columna+2] == 1)
77       vecinos_df <-rbind(vecinos_df,data.frame(x= fila, y = columna+2))
78   }

```

```

79   if(esLegal(fila,columna-2,lab)){
80     if(lab[fila,columna-2] == 1)
81       vecinos_df <- rbind(vecinos_df,data.frame(x= fila , y = columna-2)
82     )
83   }
84   return(vecinos_df)
85 }
86 # Funci n principal
87 -----
88 genera_labirinto <- function(m,n,opcion){
89 #opcion 1 = return plotss, 0 = return labirinto
90 #set.seed(2025)
91
92 # Creamos un grid de m x n
93 num_fil <- m
94 num_col <- n
95 labirinto <- creaGrid(num_fil,num_col)
96
97 # Seleccionamos la celda inicial (puede ser aleatoria tambi n)
98 #inicio <- celdaAleatoria(num_fil,num_col)
99 inicio <- c(num_fil-1,2)
100
101 # Definimos celda inicial como abierta
102 labirinto[inicio[1],inicio[2]] <- 1
103
104 # Inicializamos lista de la frontera en un dataframe
105 listaFrontera <- data.frame(x = integer(), y = integer())
106 listaFrontera <- rbind(listaFrontera,frontera(inicio[1],inicio[2],
107   labirinto))
108
109 # Empezamos el contador de iteraciones
110 k=0
111
112 # Lista para guardar las matrices en cada iteraci n (para
113   visualizaciones nicamente )
114 plotss <- list()
115
116 #Mientras no est vac a la lista
117 while( nrow(listaFrontera) != 0 ){
118   k=k+1
119   # Inicializamos lista de vecinos vac a
120   vecinos_df <- data.frame(x = integer(), y = integer())
121
122   # Seleccionamos aleatoriamente un elemento de la frontera
123   tama oLista <- nrow(listaFrontera)

```

```

122  celda <- sample(1:tama oLista , size = 1, replace = F)
123
124  # Calculamos a los vecinos del elemento frontera
125  vecinos_df <- rbind(vecinos_df, vecinos(listaFrontera[celda,1],
126                                     listaFrontera[celda,2],laberinto))
127
128  # Seleccionamos un vecino aleatoriamente y lo conectamos con la
129  # frontera
130  if(nrow(vecinos_df) != 0)
131  {
132      vecino <- sample(1:nrow(vecinos_df),size = 1)
133      laberinto[(vecinos_df[vecino,1]+listaFrontera[celda,1])/2, (
134      vecinos_df[vecino,2]+listaFrontera[celda,2])/2] = 1
135      laberinto[listaFrontera[celda,1],listaFrontera[celda,2]] <- 1
136      #laberinto[vecinos_df[vecino,1],vecinos_df[vecino,2]] <- 1
137
138  # Calculamos la frontera de la celda frontera actual y agregamos a la
139  # lista
140  listaFrontera <- rbind(listaFrontera,frontera(listaFrontera[celda,1],
141      listaFrontera[celda,2], laberinto))
142
143  # Quitamos la frontera anterior de la lista
144  listaFrontera <- listaFrontera[-celda,]
145
146  # Cuidamos que no se repitan elementos en la lista
147  listaFrontera <- unique(listaFrontera)
148  }
149
150  # Guardamos el plot (para animaci n nicamente )
151  plotss[[k]] <- laberinto
152  }
153
154  if(opcion == 1){
155
156  # png("plot1_prim.png",width = 1800, height = 1300,res=250)
157  # plot( plotss[[1]],col = c("black","white"),key = NULL,breaks = NULL
158  # ,axis.col=NULL, axis.row=NULL,xlab = "", ylab="", main = "")
159  # dev.off()
160  #
161  # png("plot2_prim.png",width = 1800, height = 1300,res=250)
162  # plot( plotss[[floor(k/2)]],col = c("black","white"),key = NULL,
163  # breaks = NULL,axis.col=NULL, axis.row=NULL,xlab = "", ylab="", main
164  # = "" )
165  # dev.off()
166  #
167  # png("plot3_prim.png",width = 1800, height = 1300,res=250)
168  # plot( plotss[[k]],col = c("black","white"),key = NULL,breaks = NULL

```

```

    ,axis.col=NULL, axis.row=NULL,xlab = "", ylab="", main = "")
161 # dev.off()
162 # En caso de querer guardar un gif de la creaci n
163 saveGIF({
164     for (i in 1:k) plot( plotss[[i]],col = c("black","white"),key =
        NULL,breaks = NULL,axis.col=NULL, axis.row=NULL,xlab = "", ylab="")
165 },interval = 0.01)
166 return(plotss)
167 }
168 else
169     return(laberinto)
170 }

```