

# Q-Learning Based Solution For Solving Mazes

Thang Doi

Data Science and Artificial Intelligence

HUST, School of Information Communication and Technology

December 29, 2024

## I. Introduction

The maze problem is a classic challenge in artificial intelligence and robotics. The primary objective is to guide an agent from a starting point to a target point within the maze. This problem can be divided into two subproblems: solving mazes in deterministic and non-deterministic environments.

In deterministic environments, the agent has complete knowledge of the maze, including all states and goals. As a result, several search algorithms, such as A\* and the Depth-First Search (DFS) algorithm, can be applied to find the optimal path. However, in non-deterministic environments, the agent cannot cover the entire maze; it can only identify its current state, making it difficult to implement a search algorithm. This scenario becomes a type of decision-making problem. Therefore, the reinforcement learning algorithm is suitable for solving this problem.

In this project, the second section shows the maze problem formula, And in the Main Work section, the main algorithm includes the maze generate algorithm, and the Q-learning algorithm to solve the maze problem. Implementing code of this project store in github

## II. Problem Formular

This section illustrates the maze problem in reinforcement learning. Structured learning in reinforcement learning can be described in the figure below.

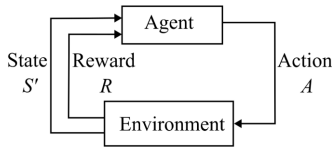


Figure 1: Reinforcement Learning Structure

Reinforcement learning [1] has several components: Agent, Environment, Action, Reward, and State. It can be represented as follows: for each state in the environment, the

agent takes an action, receives a reward from the environment, and transitions to a new state. Therefore, the learning process can be represented by a Markov Decision Process (MDP)  $(S_t, A_t, R_t, S_{t+1})$ , where the next state depends only on the current state. In the case of the maze problem, the agent starts at a random location, and the learning process terminates when the agent reaches the goal location. This project uses value-based reinforcement learning to estimate the reward of each state and action from current to future. We can see in the formula show below.

$$\begin{aligned} Q(s, a) &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^T r_{t+T} \\ Q(s, a) &= r_t + \gamma(r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{T-1} r_{t+T}) \\ Q(s, a) &= r_t + \gamma Q(s', a) \end{aligned}$$

$$\text{Temporal Different} = r_t + \gamma Q(s', a) - Q(s, a)$$

Where  $Q(s, a)$  is the expectation reward of state  $s$  when taking action  $a$ , in the Q-learning algorithm, temporal differences mainly used to update the Q-value table.

## II. Main Work

The following section shows the main work in this project, including a maze generate problem, and reinforcement learning algorithm in training agent.

### 2.1 Maze Generate Algorithm

The Maze generate problem is a classical problem in graph theory. A maze can be formally defined as a structure consisting of a set of cells, typically represented in a two-dimensional grid, where some cells are designated as walls and others as paths. The purpose of this problem breaking walls between each cell.

The maze can be represented as a graph, where the cells are vertices, and the edges represent possible paths between adjacent cells. Therefore, the solution to the maze generation problem can be modeled as a spanning tree or a cycle graph, where a spanning tree ensures connectivity without

loops, while a cycle graph includes cycles within the maze, it mean that a maze with just one path from start to goal, and otherwise more than one path. In this work, implementing two algorithm Depth First Search (DFS) and Prime algorithm for maze generation.

The first is DFS [2] with no cycle path in maze, agent can only forward one path from start to goal.

---

**Algorithm 1** DFS Maze Generation

---

```

1: Input:  $N_c$  (Number of columns),  $N_r$  (Number of rows)
2: Output: A randomly generated maze
3: Initialize  $M$  as a grid maze of size  $(N_r \times 2 - 1) \times (N_c \times 2 - 1)$  with all cells set to  $-1$ 
4: Initialize  $G$ : a graph with  $N_c \times N_r$  nodes, and  $L$ : a list to mark unvisited nodes
5:  $\text{curr\_node} = \text{start}$  as a random selection from  $G$  ( $S_{i,j} \in G$ )
6: while  $L$  is not empty do
7:    $\text{adjacents} = \text{valid\_adjacent}(\text{curr\_node})$ 
8:   if  $\text{adjacents}$  is empty then
9:     Remove  $\text{curr\_node}$  from  $L$ 
10:    if  $L$  is empty then
11:      Set  $\text{goal} = \text{curr\_node}$ 
12:      break
13:    end if
14:    continue
15:  else
16:     $\text{next\_node} = \text{random\_selection}(\text{adjacents})$ 
17:     $G[\text{next\_node}] = \text{curr\_node}$ 
18:     $\text{curr\_node} = \text{next\_node}$ 
19:  end if
20: end while
21: // Using the spanning tree to break walls in the grid maze
22: // Node in graph mapping to maze coordinates:  $(i, j) \rightarrow (i \times 2, j \times 2)$ 
23: for each node in  $G$  do
24:    $M[i * 2][j * 2] = 0$   $\triangleright$  Mark the node as open in the maze
25:    $\text{node\_connect} = G[\text{node}]$ 
26:    $\text{dx}, \text{dy} = \text{distance}(\text{node}, \text{node\_connect})$ 
27:    $M[\text{node}[0] + \text{dx}][\text{node}[1] + \text{dy}] = 0$   $\triangleright$  Break the wall between nodes
28: end for

```

---

The second algorithm is Prim's [2] Algorithm, used to generate a spanning tree. The process begins by creating a spanning tree based on Prim's method, which ensures that all nodes are connected with no cycles. After constructing the spanning tree, random reconnections are introduced to create cycles within the graph. This means that the agent can have multiple paths from the start node to the goal node, offering

greater flexibility and variability in navigating the maze.

---

**Algorithm 2** Prime Maze Cycle Generation

---

**Require:**  $N_c$  (Number of columns),  $N_r$  (Number of rows)

**Ensure:** A randomly generated cycle maze

```

1: Initialize  $M$  as a grid maze of size  $(N_r \times 2 - 1) \times (N_c \times 2 - 1)$  with all cells set to  $-1$ 
2: Initialize  $G$ : a graph with  $N_c \times N_r$  nodes
3: Initialize  $L$ : a list to mark connected edges, and  $N_e$ : number of edges created, set to 0
4: while  $N_e < N_c \times N_r$  do
5:    $\text{node} \leftarrow \text{random\_node}(G)$ 
6:    $\text{list\_adjacents} \leftarrow$  adjacent nodes of node in  $L$ 
7:   if  $\text{list\_adjacents}$  is empty then
8:     remove node from  $G$ 
9:     continue
10:  else
11:     $\text{adjacent} \leftarrow \text{random\_node}(\text{list\_adjacents})$ 
12:  end if
13:  if not  $\text{is\_created\_cycle}(\text{node})$  then
14:     $G[\text{node}][\text{adjacent}] \leftarrow \text{False}$ 
15:     $L[\text{node}][\text{adjacent}] \leftarrow \text{True}$ 
16:    Increment  $N_e$ 
17:  end if
18: end while
19:  $\text{create\_cycle}(G, \text{number\_cycle})$ 
20: Perform wall-breaking process (like DFS)
21: return  $M$ 

```

---

## 2.2 Q-Learning Algorithm

Q-learning [3] is a fundamental reinforcement learning algorithm used to train agents to make sequential decisions in an environment. It is a model-free method, meaning it does not require prior knowledge of the environment's dynamics, and it seeks to find an optimal policy that maximizes the cumulative reward over time.

---

**Algorithm 3** Q-Learning Algorithm

---

```
1: Input: State space  $S$ , Action space  $A$ , Learning rate  $\alpha$ ,  
   Discount factor  $\gamma$ , Exploration rate  $\epsilon$ , Number of episodes  
    $N$   
2: Output: Q-table  $Q(s, a)$   
3: Initialize  $Q(s, a) =$  arbitrarily for all  $s \in S, a \in A$   
4: for episode = 1 to  $N$  do  
5:   current state  $s =$  start location  
6:   while  $s$  is not terminal do  
7:     if  $s$  not in  $Q(s, a)$ , then update  $(s, valid_a)$   
8:      $t$  : random probability  
9:      $a = \begin{cases} \text{a random action } a \in A_s, & \text{if } t > \epsilon, \\ \arg \max_a Q(s, a), & \text{otherwise.} \end{cases}$   
10:    Take action  $a$ , observe reward  $r$  and next state  $s'$   
    from environment  
11:    Update  $Q(s, a)$ :  

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$
  
12:     $s \leftarrow s'$   
13:   end while  
14: end for  
15: return  $Q(s, a)$ 
```

---

Double Q-learning [4] is an enhancement of the traditional Q-learning algorithm designed to address the overestimation bias that can occur when updating Q-values. In standard Q-learning, the maximum Q-value  $\max_a Q(s, a)$  is used to estimate the future reward, which may lead to overly optimistic value estimates.

Double Q-learning mitigates this issue by maintaining two separate Q-tables  $Q_1$  and  $Q_2$ . During updates, one Q-table is used to select the action, and the other is used to evaluate it.

---

**Algorithm 4** Double Q-Learning Algorithm

---

```
1: Input: State space  $S$ , Action space  $A$ , Learning rate  $\alpha$ ,  
   Discount factor  $\gamma$ , Exploration rate  $\epsilon$ , Number of episodes  
    $N$   
2: Output: Two Q-table  $Q_1(s, a), Q_2(s, a)$   
3: Initialize  $Q_1(s, a) = , Q_2(s, a) =$  arbitrarily for all  $s \in S,$   
    $a \in A$   
4: for episode = 1 to  $N$  do  
5:   current state  $s =$  start location  
6:   while  $s$  is not terminal do  
7:     if  $s$  not in  $Q_1(s, a), Q_2(s, a)$  then update  $(s, valid_a)$   
8:      $t$  : random probability  
9:      $a = \begin{cases} \text{a random action } a \in A_s, & \text{if } t > \epsilon, \\ \arg \max_a Q_1(s, a), & \text{if } table - select = 0, \\ \arg \max_a Q_2(s, a), & \text{if } table - select = 1 \end{cases}$   
10:    Take action  $a$ , observe reward  $r$  and next state  $s'$   
    from environment  
11:    If table-select == 0 :  

$$Q_1(s, a) \leftarrow Q_1(s, a) + \alpha [r + \gamma \max_{a'} Q_2(s', a') - Q_1(s, a)]$$
  
12:    If table-select == 1  

$$Q_2(s, a) \leftarrow Q_2(s, a) + \alpha [r + \gamma \max_{a'} Q_1(s', a') - Q_2(s, a)]$$
  
13:     $s \leftarrow s'$   
14:   end while  
15: end for  
16: return  $Q_1(s, a), Q_2(s, a)$ 
```

---

### III. Results and Application

To visualize the agent's learning process in a maze environment, the project implements a GUI main program. The structure and functionality of the program are illustrated in the figure below.

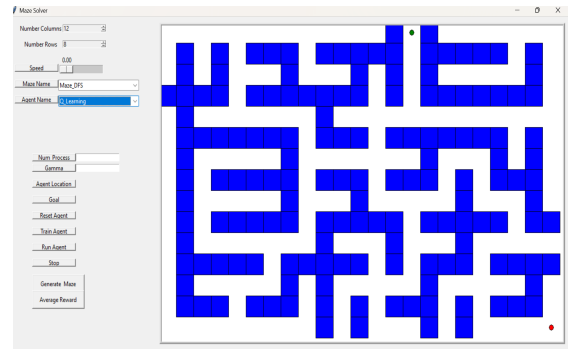


Figure 2: GUI project

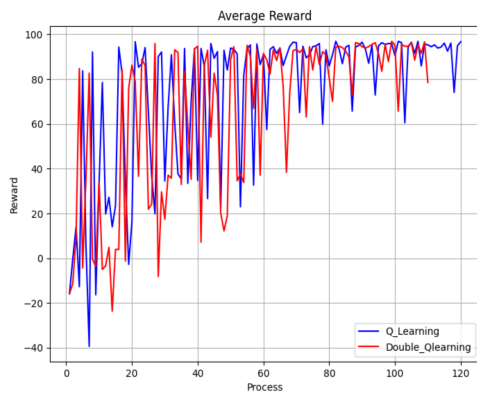


Figure 3: Average reward of Q-learning and Double Q-learning Algorithm

## 1 IV. Conclusion

This project implement a simple reinforcement learning algorithm that is my interest topic. Q-learning is a tabular algorithm, therefore it have several limitation such as inefficient in learning with large state and action space. In the future, I decide develop enhancent algorithm deep reinforcement learning having SOTA compare with traditional reinforcement.

## Acknowledgments

I would like to express our sincere gratitude to Teacher Michel Toulouse, my instructor for the final semester Project 1, for providing me with excellent guidance and support, enabling us to successfully propose and complete this project.

## References

- [1] Sutton, R. S., & Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. The MIT Press.
- [2] Wikipedia contributors. (n.d.). *Maze generation algorithm*. Retrieved from <https://en.wikipedia.org/wiki/Maze-generation-algorithm>
- [3] Watkins, C. J. C. H. (1989). *Learning from Delayed Rewards*. PhD thesis, University of Cambridge.
- [4] Van Hasselt, H. (2010). *Double Q-learning*. In *Proceedings of the 23rd International Conference on Neural Information Processing Systems (NIPS 2010)*.