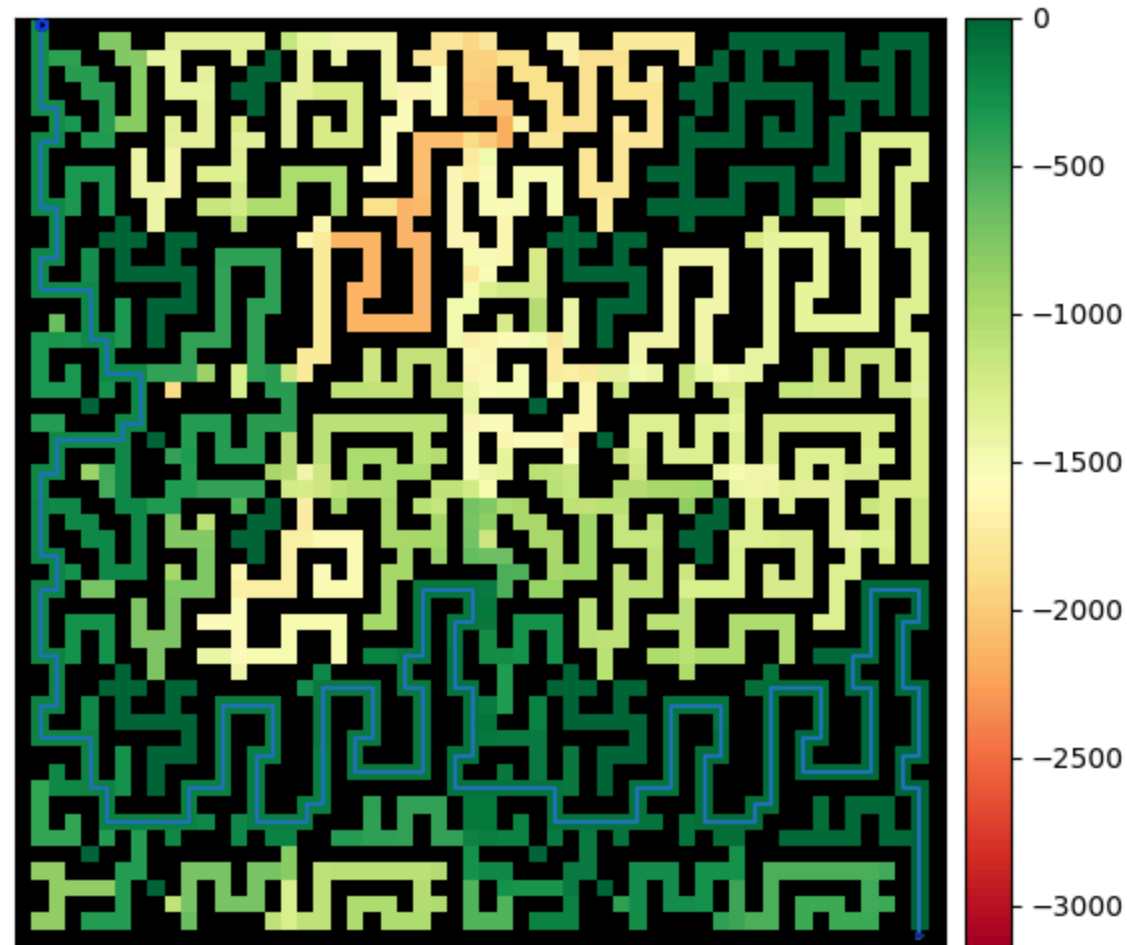


Solving a Maze - Reinforcement Learning



Contents

- ▶ Reinforcement Learning - Introduction
- ▶ Environment: Maze
- ▶ Agent: Robot
- ▶ Reward System
- ▶ Results and Improvements

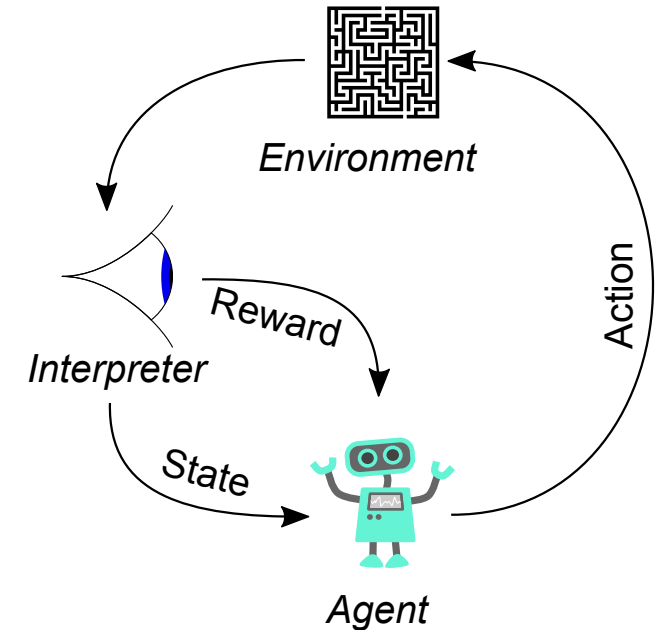
Reinforcement Learning

(Un)supervised

- ▶ Need of a training set
 - ▶ May be labeled (supervised)
 - ▶ Or unlabeled (unsupervised)

Reinforcement

[...] intelligent agent ought to take actions in a dynamic environment in order to maximize the cumulative reward.
-- Wikipedia



The Environment

- ▶ Provide a system, with which the agent can interact
- ▶ Implement game logic

```
class Maze():  
  
    def get_moves():  
    def move(...):  
  
    def get_state_and_reward():  
  
    def is_game_over():  
  
    # just for plotting  
    def get_steps()
```

- ▶ Stores:
 - ▶ Actual maze
 - ▶ Player position
- ▶ Provides to agent:
 - ▶ Possible moves
 - ▶ Rewards
- ▶ Ends game

The Agent

- ▶ Robot which learns to solve the maze

```
class Agent():  
    def __init__(...):  
        """set learning/exploration rate"""  
  
    def choose_action(...):  
  
    def store_rewards(...):  
  
    def learn():
```

- ▶ Initialization with specific
 - ▶ learning rate
 - ▶ exploration rate
- ▶ Choose from possible moves
- ▶ Store reward for given position
- ▶ update reward table from rewards collected during one game (learn)

Reward System

Environment

```
def get_state_and_reward(self):  
    return self.pos, -1
```

- ▶ penalize each step with -1

Agent

- ▶ During the game

```
def store_reward(self, pos, reward):  
    self._pos_history.append(tuple(pos))  
    self._rew_history.append(reward)
```

- ▶ After the game `learn()`

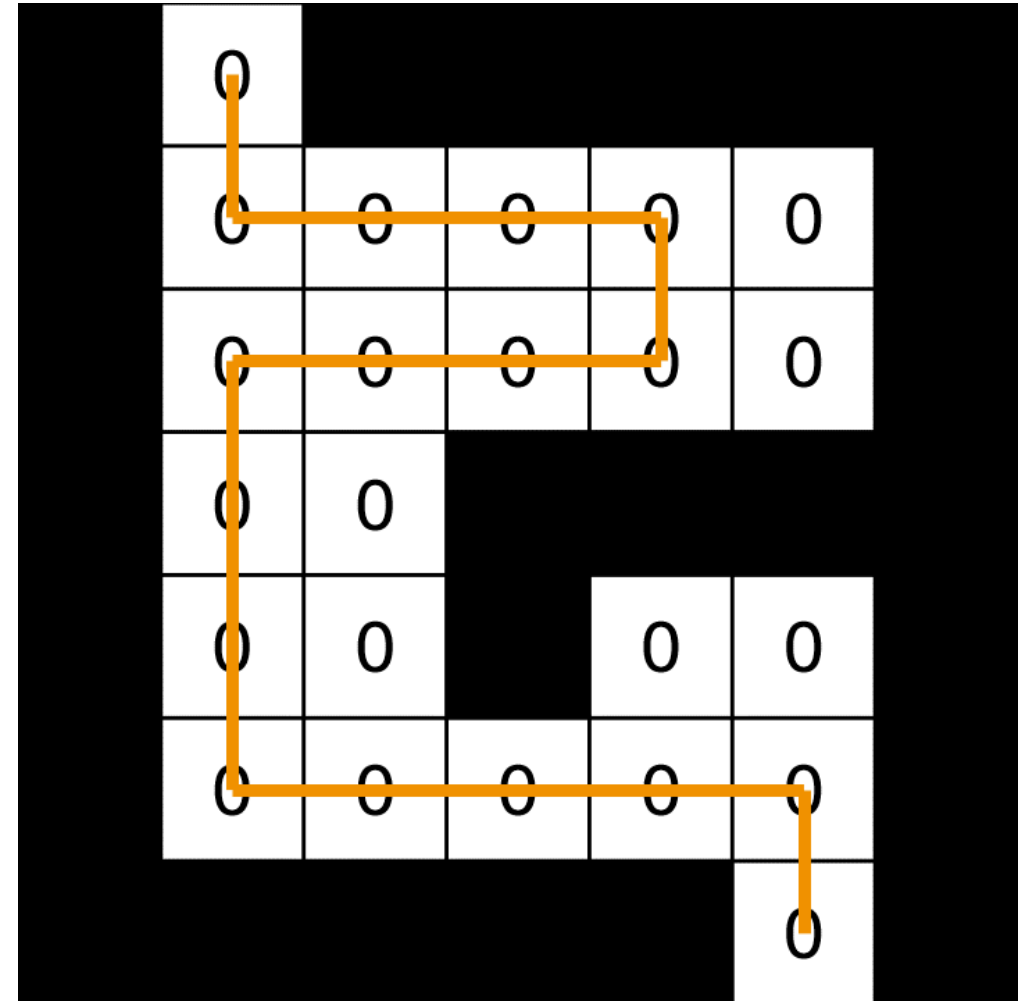
Reward System - *Agent*

- ▶ Agent stores a reward matrix
 - ▶ technically **G table** (*cumulative*)
- ▶ initialized to zero

0				
0	0	0	0	0
0	0	0	0	0
0	0			
0	0		0	0
0	0	0	0	0
				0

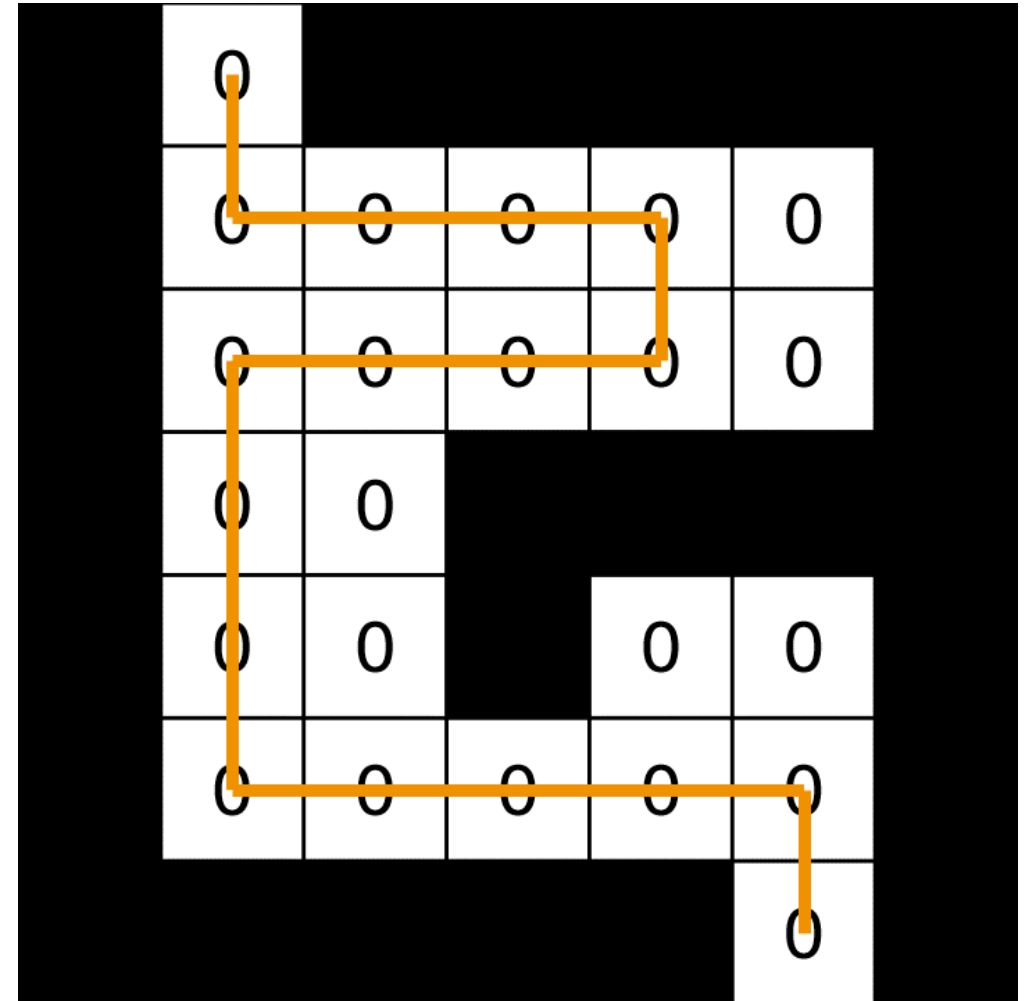
Reward System - *Agent*

- ▶ First game is random



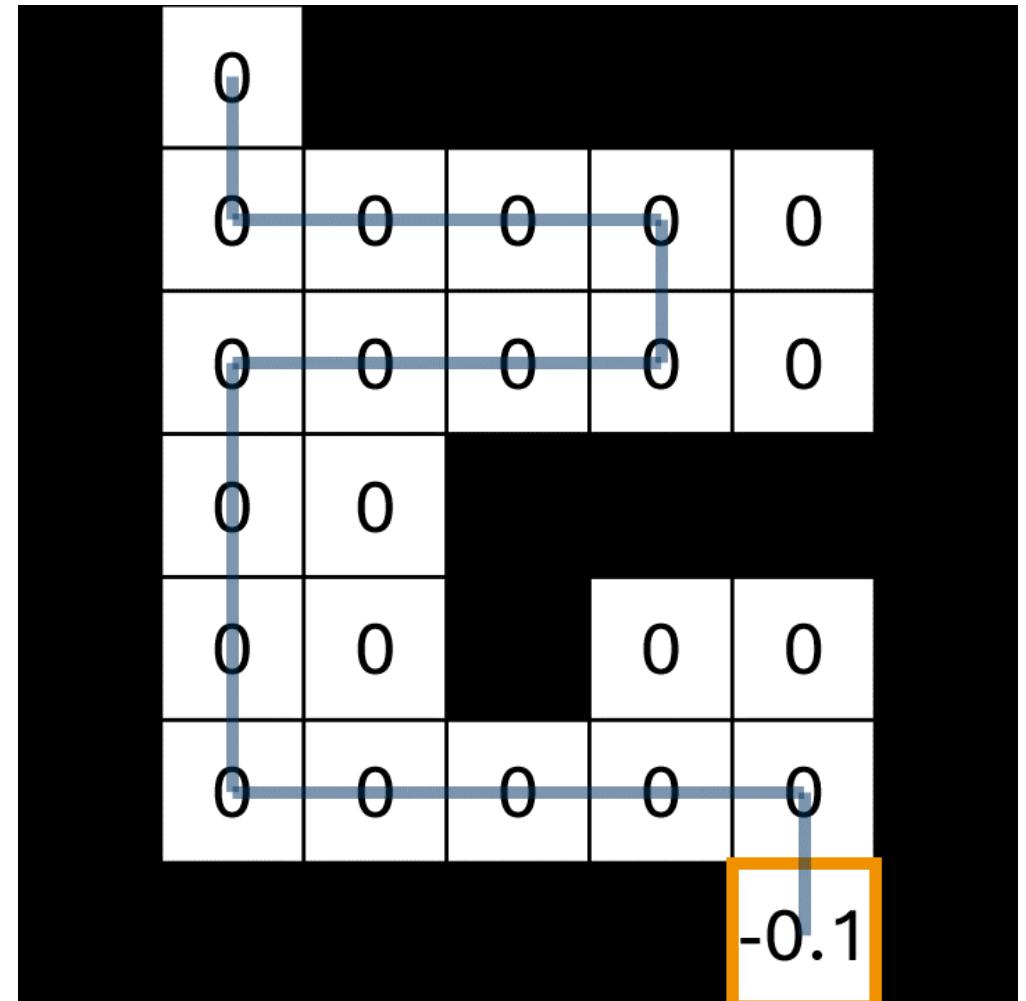
Reward System - *Agent - learning*

$$G_{\text{path}} = 0$$



Reward System - *Agent - learning*

- Loop path **backwards**



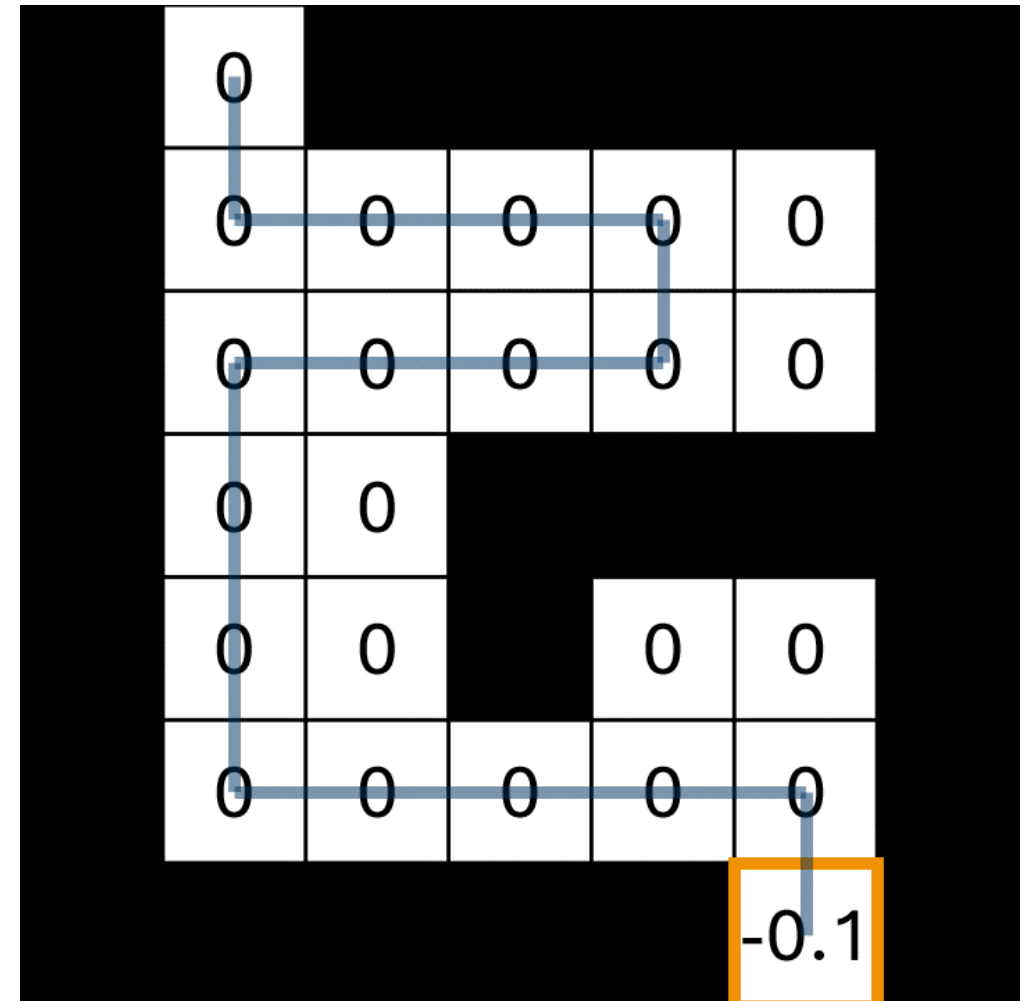
Reward System - *Agent - learning*

- ▶ Loop path **backwards**
- ▶ Compute cumulative reward along path

$$reward_{pos} = -1$$

$$G_{path} = G_{path} + reward_{pos}$$

$$G_{path} = -1$$



Reward System - *Agent - learning*

- ▶ Loop path **backwards**
- ▶ Compute cumulative reward along path

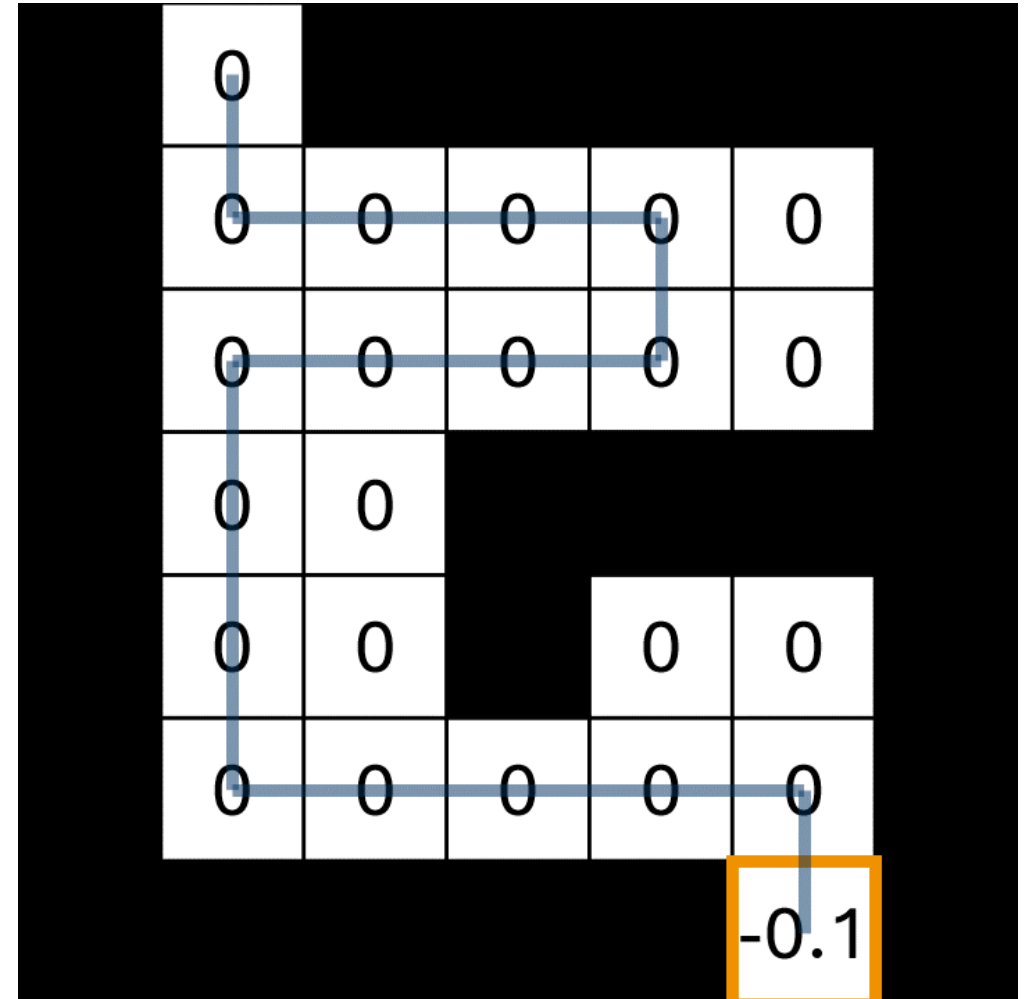
$$reward_{pos} = -1$$

$$G_{path} = G_{path} + reward_{pos}$$

$$G_{path} = -1$$

- ▶ Update reward matrix entry (eg. $\alpha = 0.1$)

$$G_{pos} = G_{pos} + \alpha(G_{path} - G_{pos}) = -0.1$$



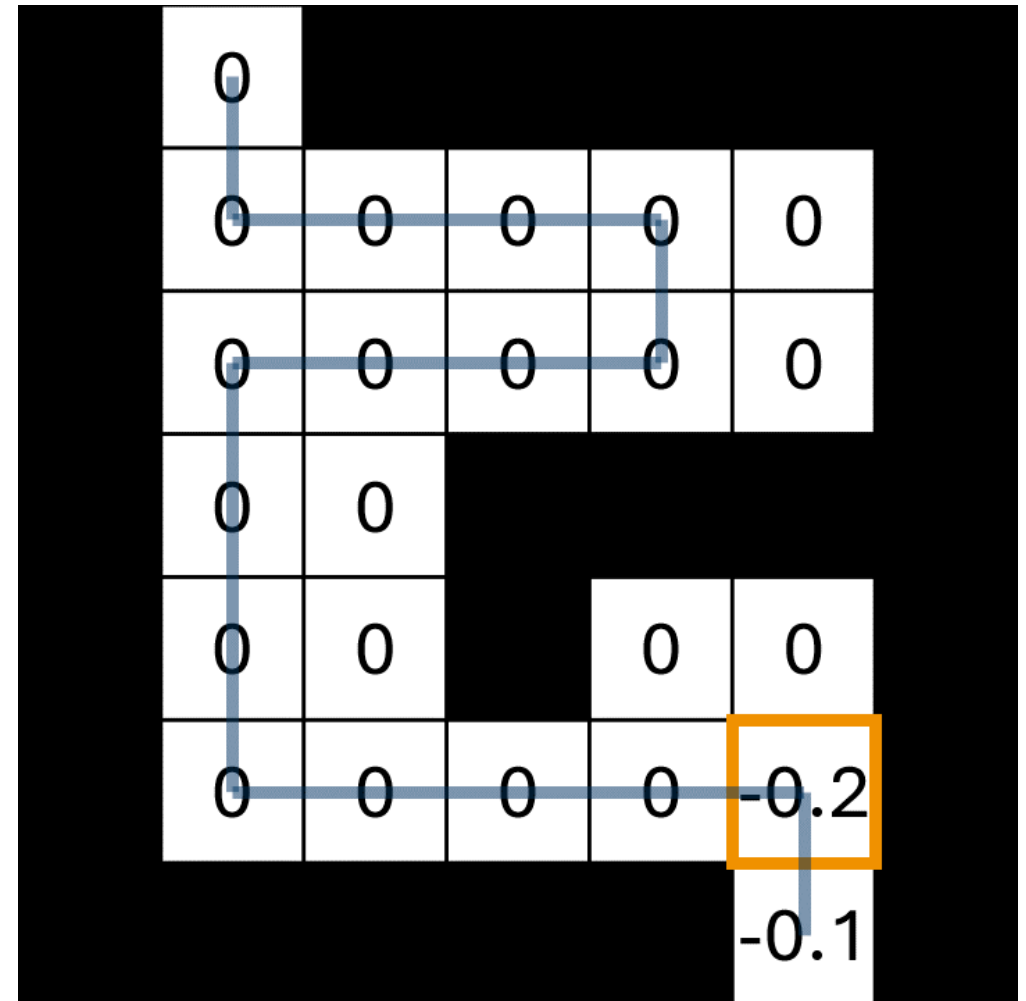
Reward System - *Agent - learning*

$$reward_{pos} = -1$$

$$G_{path} = G_{path} + reward_{pos}$$

$$G_{path} = -2$$

$$G_{pos} = G_{pos} + \alpha(G_{path} - G_{pos}) = -0.2$$



A 7x7 grid representing a sparse matrix. The matrix is mostly zeros, with a few non-zero values: -0.3 at (row 6, column 5) and -0.2 at (row 6, column 6). The cell containing -0.3 is highlighted with an orange square. Blue lines connect the zeros in a path from (1,1) to (6,6).

0						
0	0	0	0	0	0	
0	0	0	0	0	0	
0	0					
0	0			0	0	
0	0	0	-0.3	-0.2		
					-0.1	

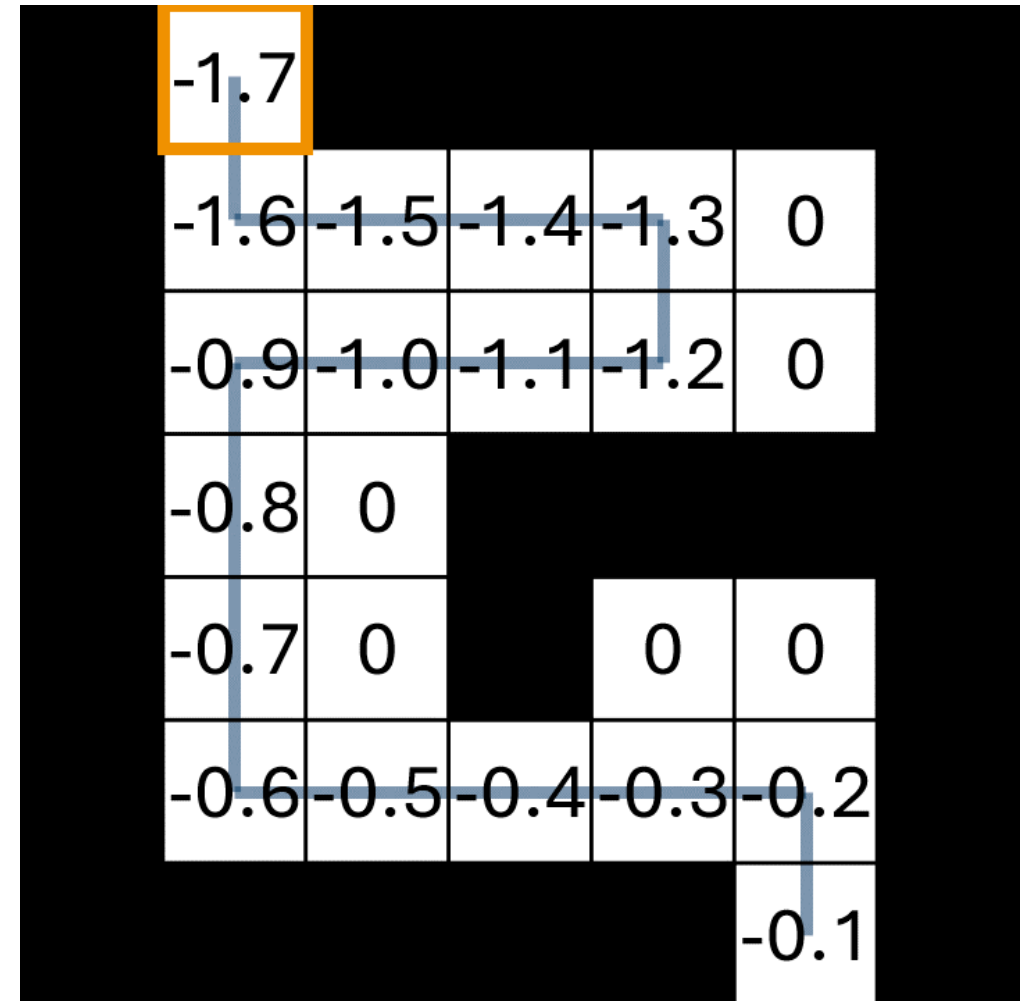
Reward System - *Agent - learning*

$$reward_{pos} = -1$$

$$G_{path} = G_{path} + reward_{pos}$$

$$G_{path} = -17$$

$$G_{pos} = G_{pos} + \alpha(G_{path} - G_{pos}) = -1.7$$



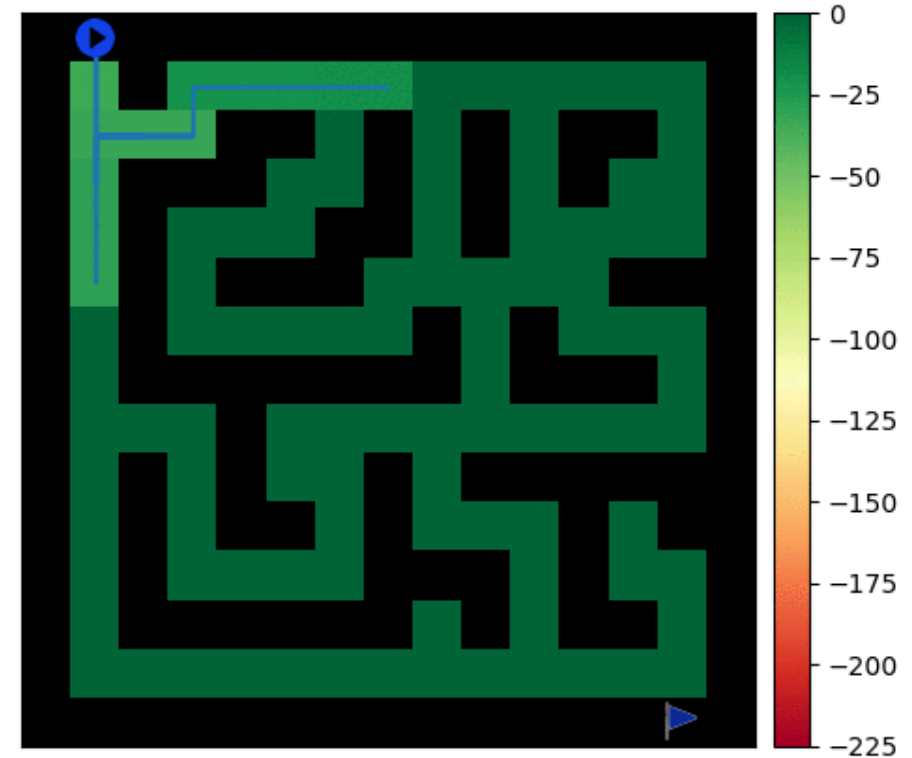
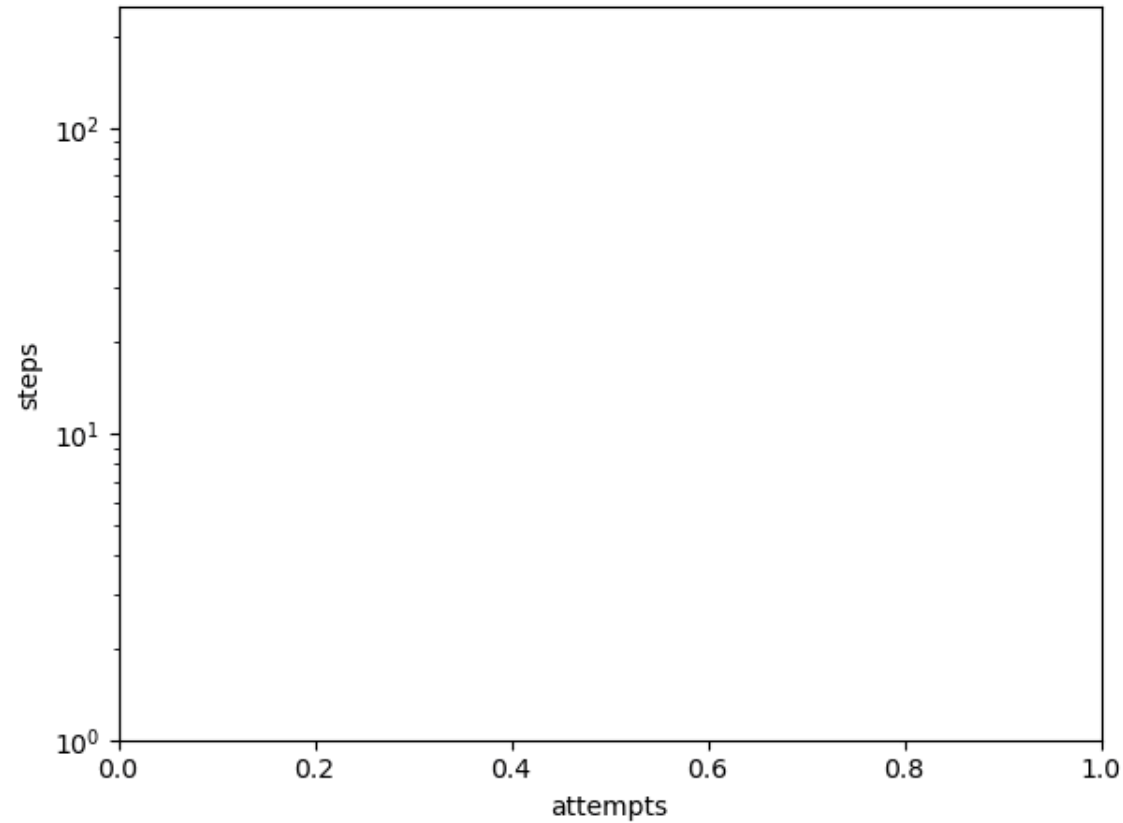
Reward System - *Agent - choosing*

- ▶ Consider second run
- ▶ Agent has learned (some) reward matrix
- ▶ Choose action towards maximum reward
- ▶ Exploration rate:
 - ▶ Choose random action by chance

-1.7				
-1.6	-1.5	-1.4	-1.3	0
-0.9	-1.0	-1.1	-1.2	0
-0.8	0			
-0.7	0		0	0
-0.6	-0.5	-0.4	-0.3	-0.2
				-0.1

Small Maze

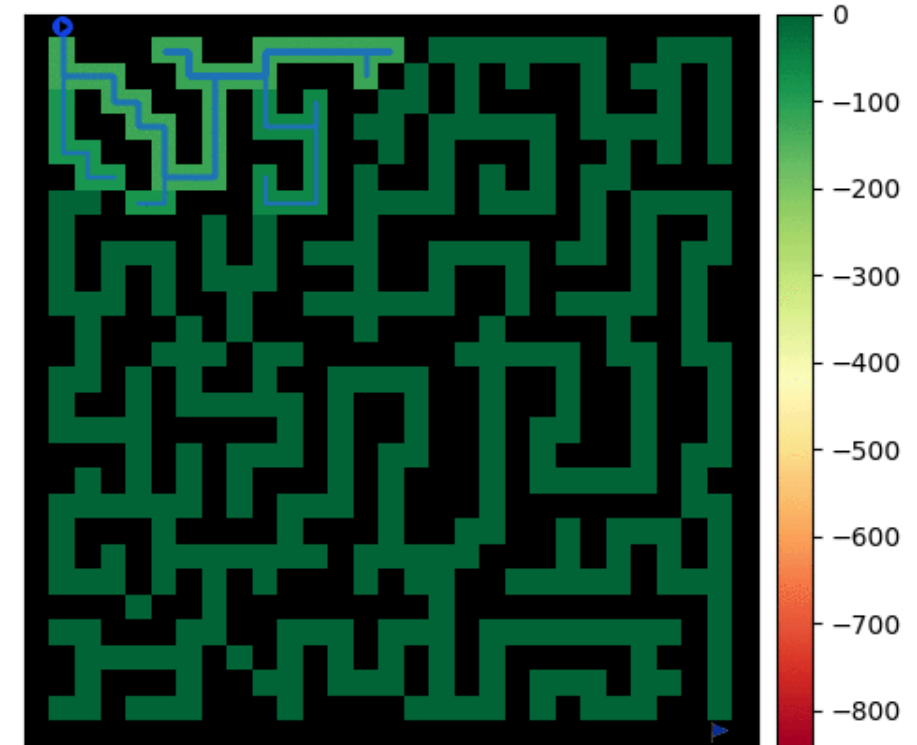
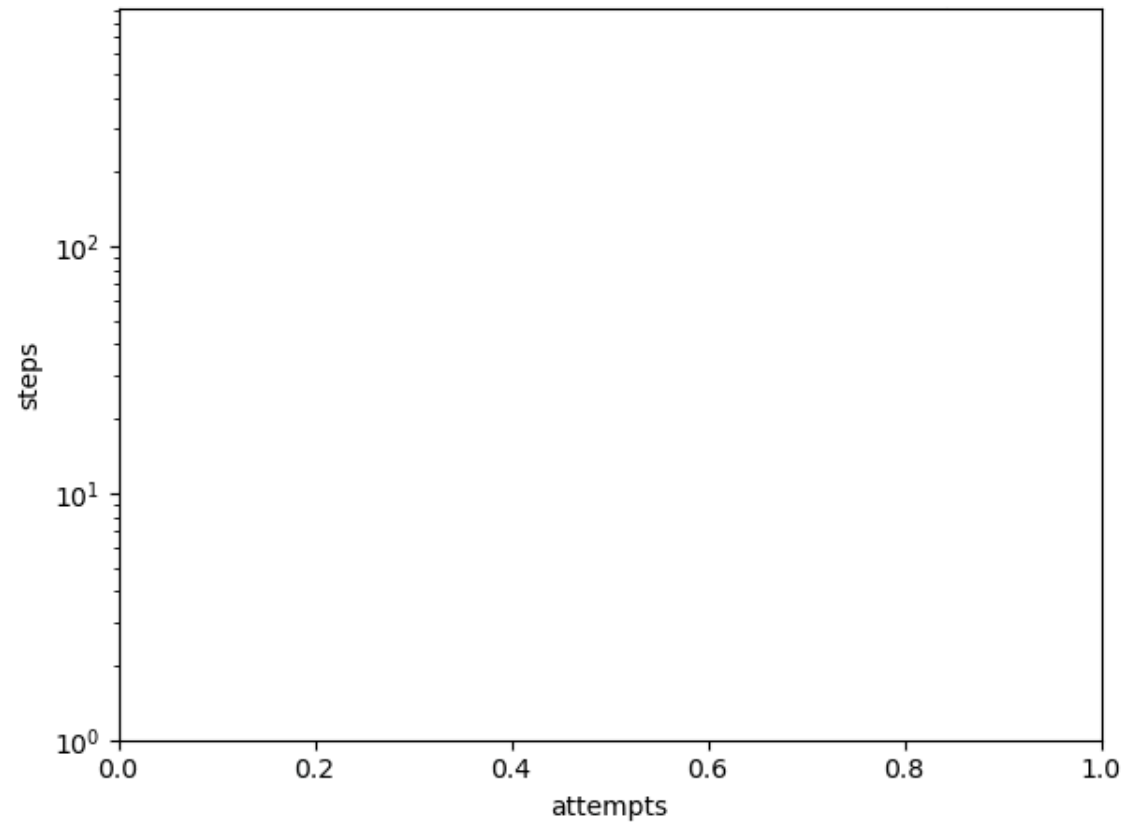
Learning Rate: 0.15 Exploration Rate: 0.25



► Remove Backtracking

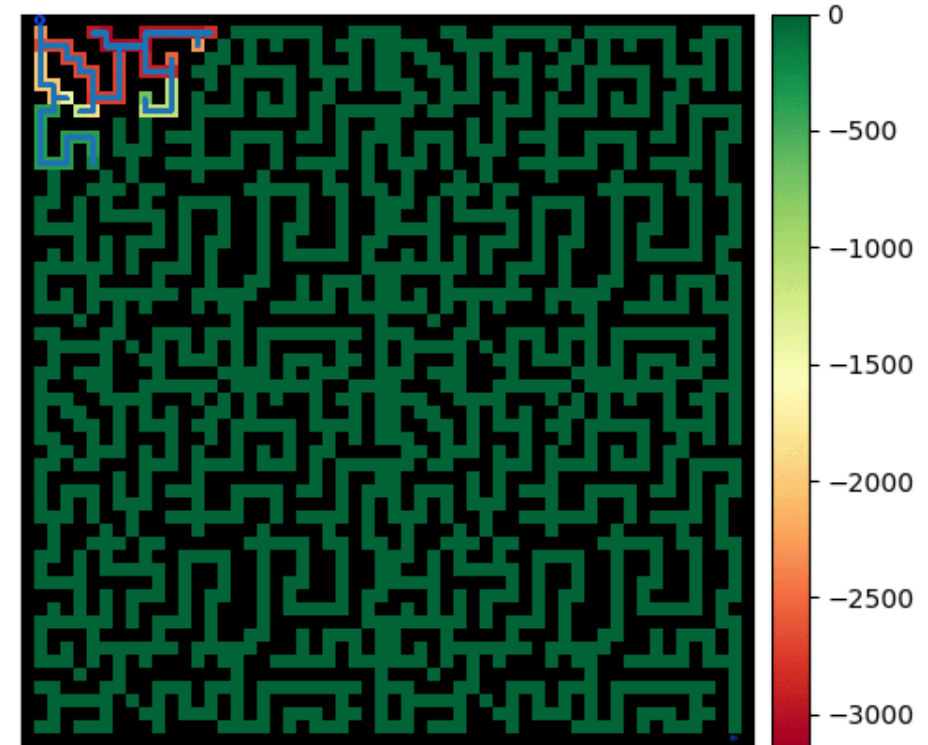
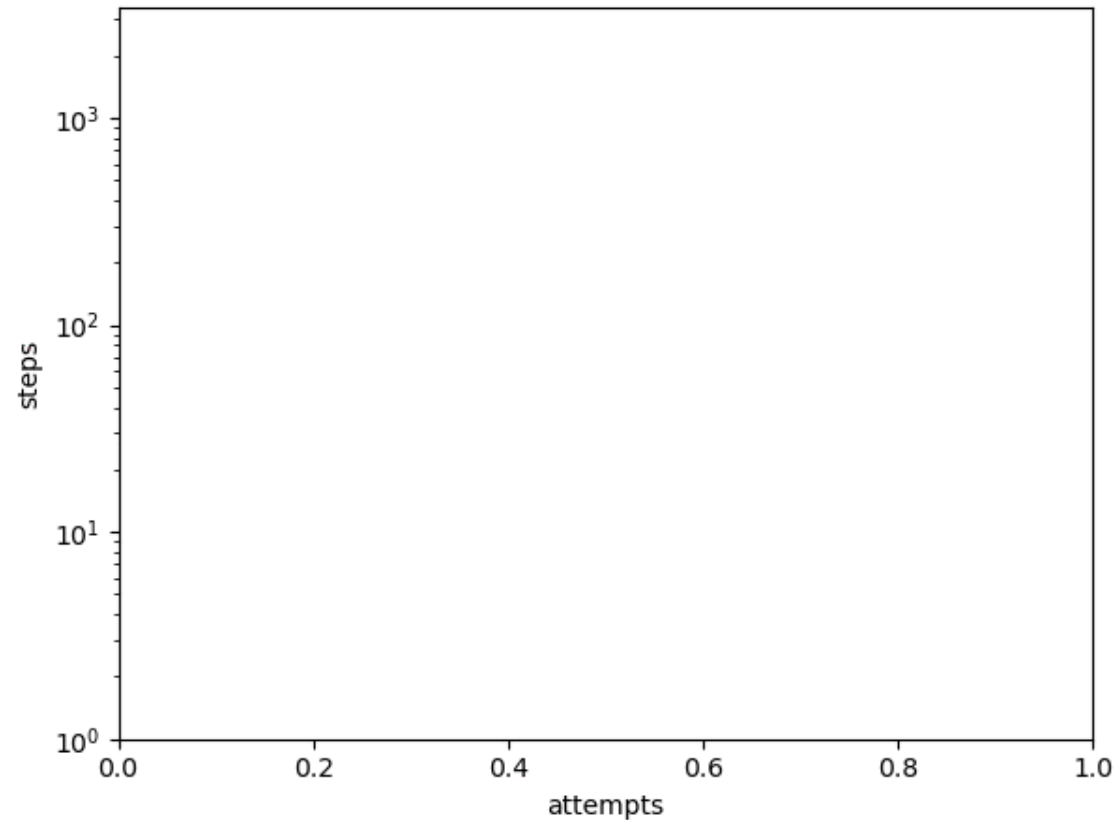
Bigger!

Learning Rate: 0.15 Exploration Rate: 0.25



And BIGGER!

Learning Rate: 0.15 Exploration Rate: 0.25



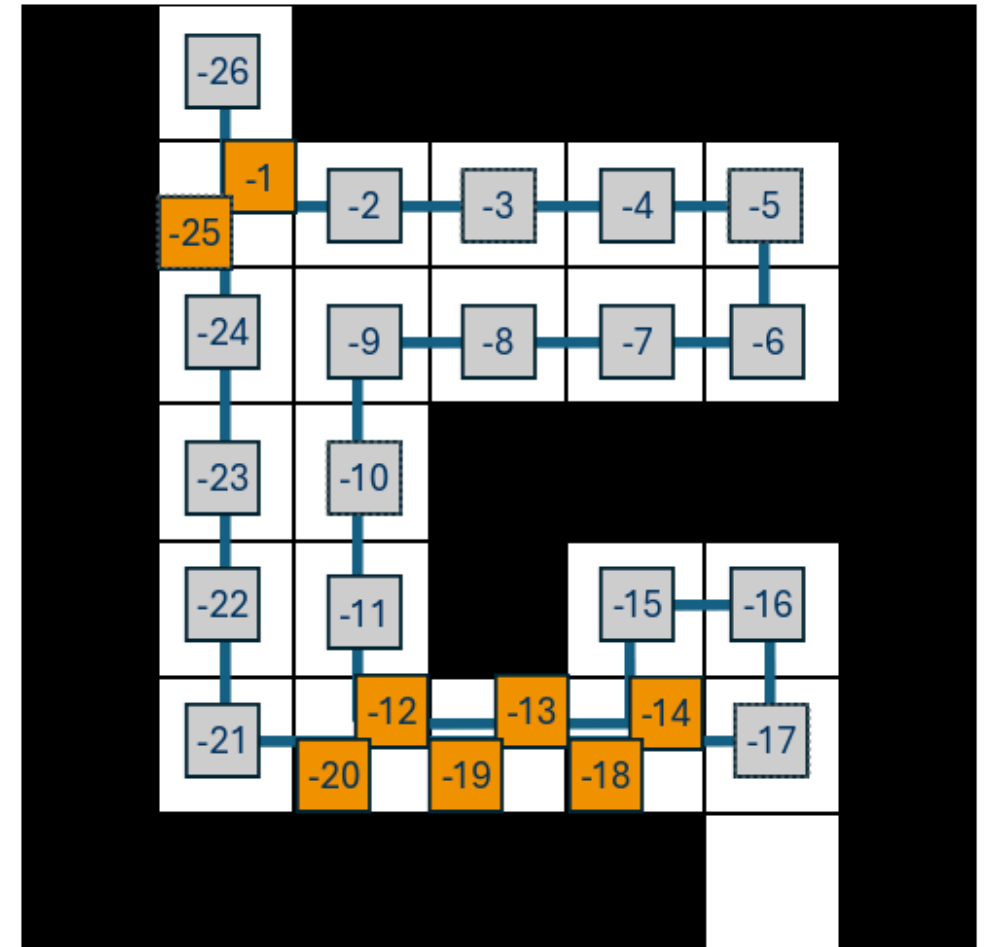
- Why aren't dead ends deep red?

Cumulative Reward

$$G_{\text{pos}} = G_{\text{nos}} + \alpha(G_{\text{path}} - G_{\text{pos}})$$

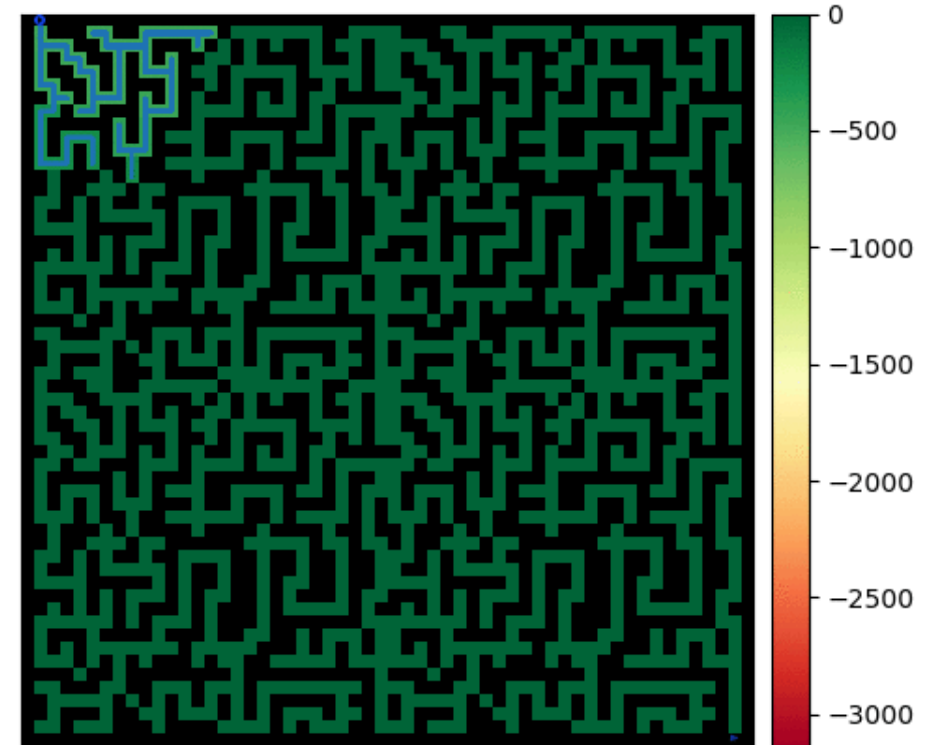
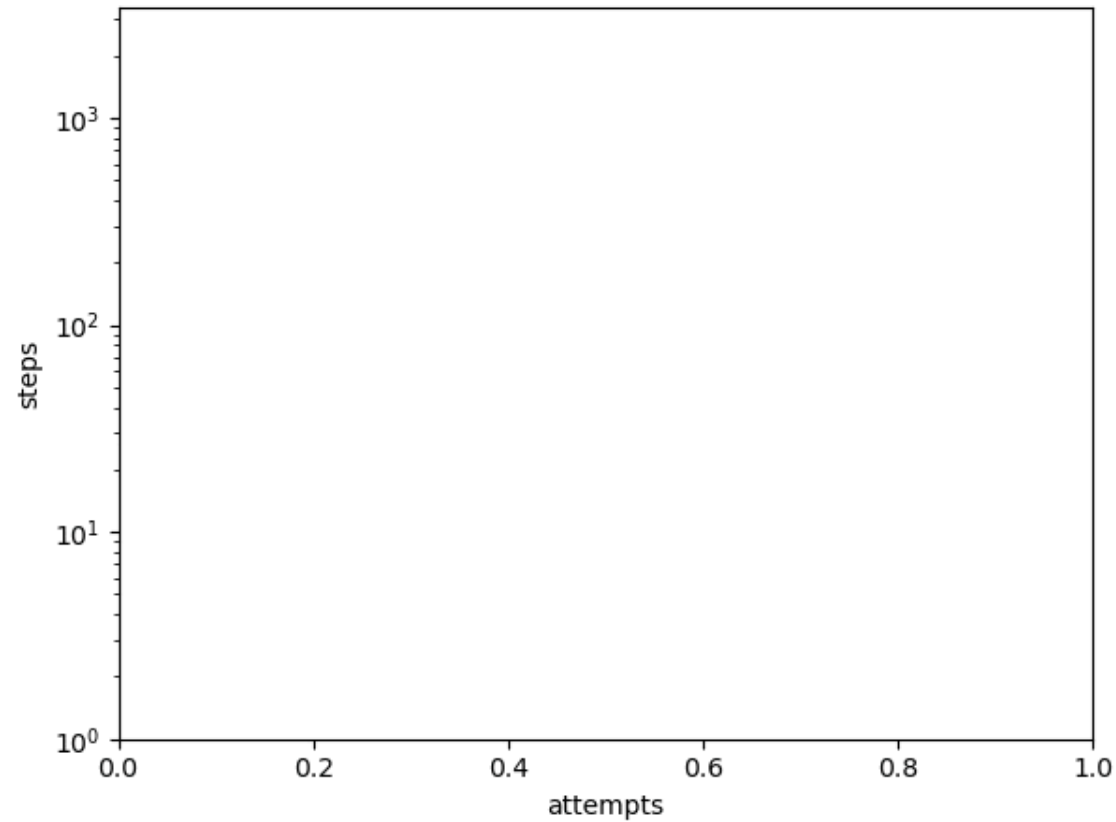
\Rightarrow

$$G_{\text{pos}} = G_{\text{pos}} + \alpha(\text{min}\{G_{\text{path}}\} - G_{\text{pos}})$$



Hurray!

Learning Rate: 0.15 Exploration Rate: 0.25



Thanks for your attention

Appendix

Removing Backtracking

- ▶ If action returns to previous position
 - ▶ Remove action from `action_list` and choose again (except for dead end)

```
def chose_action(action_list):  
    ...  
    # attempt to make him not reverse moves immediatly  
    # If action returns to previous position  
    # - remove action from action list if possible (ie. not only option)  
    # - choose again  
    if tuple(action_list[index]()) == self._pos_history[-2]:  
        if len(action_list) > 1:  
            action_list.pop(index)  
        return self.choose_action(action_list)
```

- ▶ May not backtrack directly
- ▶ Can still run in loops