

人工智能与机器学习报告

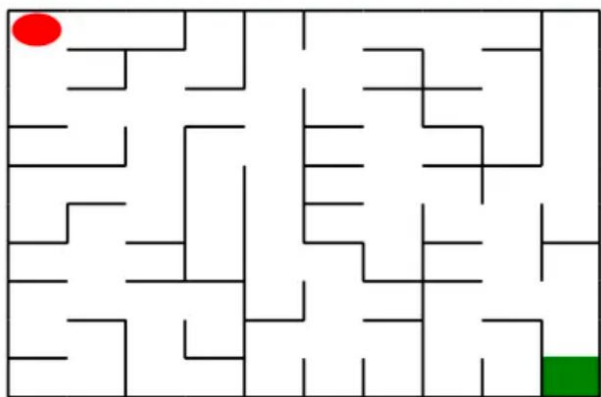
2. 机器人走迷宫

3220101111 洪晨辉

1. 实验介绍

1.1 实验背景

本实验分别使用基础搜索算法和 Deep QLearning 算法，完成机器人自动走迷宫。



如上图所示，左上角的红色椭圆既是起点也是机器人的初始位置，右下角的绿色方块是出口。

游戏规则为：从起点开始，通过错综复杂的迷宫，到达目标点(出口)。在任一位置可执行动作包括：向上走 'u'、向右走 'r'、向下走 'd'、向左走 'l'。执行不同的动作后，根据不同的情况会获得不同的奖励。

1.2 实验要求

- A)使用 Python 语言。
- B)使用基础搜索算法完成机器人走迷宫。
- C)使用 Deep QLearning 算法完成机器人走迷宫。
- D)算法部分需要自己实现，不能使用现成的包、工具或者接口。

1.3 实验环境

可以使用 Python 实现基础算法的实现，使用 Keras、PyTorch 等框架实现 Deep QLearning 算法。本例中采用 PyTorch。

2. 实验操作

2.1 创建迷宫

通过迷宫类 Maze 可以随机创建一个迷宫。具体实现存在 Maze.py 中，已经写好。使用 Maze(maze_size=size) 来随机生成一个定大小的迷

宫。

raw_maze(): 画出当前的迷宫

2.2 创建机器人

同理，有一个机器人类，存在 Maze.py 中。

Maze 类中重要的成员方法如下：

sense_robot()：获取机器人在迷宫中目前的位置。

return: 机器人在迷宫中目前的位置。

move_robot(direction)：根据输入方向移动默认机器人，若方向不合法则返回错误信息。

direction: 移动方向，如:"u"，合法值为：['u', 'r', 'd', 'l']

return: 执行动作的奖励值

can_move_actions(position): 获取当前机器人可以移动的方向

position: 迷宫中任一处的坐标点

return: 该点可执行的动作，如：['u', 'r', 'd']

is_hit_wall(self, location, direction): 判断该移动方向是否撞墙

location, direction: 当前位置和要移动的方向，如(0,0), "u"

return: True(撞墙) / False(不撞墙)

2.3 基础搜索算法

对于迷宫游戏，常见的三种搜索算法有广度优先搜索、深度优先搜索和最佳优先搜索（A*）。这里采用最佳优先搜索算法。当然，经过测试，这三种算法都可以很好地完成搜索的任务。

```
1. # 导入相关包
2. import os
3. import random
4. import numpy as np
5. from Maze import Maze
6. from Runner import Runner
7. from QRobot import QRobot
8. from ReplayDataSet import ReplayDataSet
```

```

9. from torch_py.MinDQNRobot import MinDQNRobot
   t as TorchRobot # PyTorch 版本
10. from keras_py.MinDQNRobot import MinDQNRobot
   ot as KerasRobot # Keras 版本
11. import matplotlib.pyplot as plt
12. from functools import total_ordering
13.
14. import numpy as np
15. import heapq
16.
17. # 机器人移动方向
18. move_map = {
19.     'u': (-1, 0), # up
20.     'r': (0, +1), # right
21.     'd': (+1, 0), # down
22.     'l': (0, -1), # left
23. }
24.
25. @total_ordering
26. class SearchTree(object):
27.     def __init__(self, loc=(), action='',
28.                  parent=None, g=0, h=0):
29.         self.loc = loc # 当前节点位置
30.         self.to_this_action = action # 到达当前节点的动作
31.         self.parent = parent # 当前节点的父节点
32.         self.children = [] # 当前节点的子节点
33.         self.g = g # 实际代价
34.         self.h = h # 启发式代价
35.
36.     def add_child(self, child):
37.         self.children.append(child)
38.
39.     def is_leaf(self):
40.         return len(self.children) == 0
41.
42.     def f(self):
43.         return self.g + self.h
44.
45.     def __eq__(self, other):
46.         return self.f() == other.f()

```

```

47.     def __lt__(self, other):
48.         return self.f() < other.f()
49.
50. def heuristic(loc, goal):
51.     """
52.     计算启发式代价函数 (曼哈顿距离)
53.     :param loc: 当前节点位置
54.     :param goal: 目标节点位置
55.     :return: 曼哈顿距离
56.     """
57.     return abs(loc[0] - goal[0]) + abs(loc
58. [1] - goal[1])
59.
60. def expand(maze, is_visit_m, node, goal):
61.     """
62.     拓展叶子节点，即为当前的叶子节点添加执行合法动作后到达的子节点
63.     :param maze: 迷宫对象
64.     :param is_visit_m: 记录迷宫每个位置是否访问的矩阵
65.     :param node: 待拓展的叶子节点
66.     :param goal: 目标位置
67.     """
68.     can_move = maze.can_move_actions(node.loc)
69.     for a in can_move:
70.         new_loc = tuple(node.loc[i] + move_map[a][i] for i in range(2))
71.         if not is_visit_m[new_loc]:
72.             g = node.g + 1 # 每移动一步的代价为1
73.             h = heuristic(new_loc, goal)
74.             child = SearchTree(loc=new_loc, action=a, parent=node, g=g, h=h)
75.             node.add_child(child)
76.
77. def back_propagation(node):
78.     """
79.     回溯并记录节点路径
80.     :param node: 待回溯节点
81.     :return: 回溯路径
82.     """
83.     path = []
84.     while node.parent is not None:

```

```

84.         path.insert(0, node.to_this_action
85.     )
86.     node = node.parent
87.     return path
88. def a_star_search(maze):
89.     """
90.     使用 A*算法对迷宫进行路径搜索
91.     :param maze: 待搜索的 maze 对象
92.     :return: 从起点到目标点的路径
93.     """
94.     start = maze.sense_robot()
95.     goal = maze.destination
96.     root = SearchTree(loc=start, g=0, h=heuristic(start, goal))
97.
98.     open_list = [] # 优先队列存储待扩展节点
99.     heapq.heappush(open_list, (root.f(), root))
100.
101.     h, w, _ = maze.maze_data.shape
102.     is_visit_m = np.zeros((h, w), dtype=np.int) # 标记迷宫的各个位置是否被访问过
103.
104.     while open_list:
105.         _, current_node = heapq.heappop(open_list)
106.         is_visit_m[current_node.loc] = 1
107.         # 标记当前节点位置已访问
108.         if current_node.loc == goal: # 到达目标点
109.             return back_propagation(current_node)
110.
111.         if current_node.is_leaf():
112.             expand(maze, is_visit_m, current_node, goal)
113.
114.         for child in current_node.children:
115.             if not is_visit_m[child.loc]:
116.                 heapq.heappush(open_list,

```

```

117.
118.         return [] # 未找到路径

```

这里，其启发式函数 $h(n)$ 为普通的曼哈顿距离。算法会从 `open_list` 中取出 $f(n)$ 最小的节点。如果当前节点位置等于目标位置，回溯输出路径。否则，扩展当前节点的子节点，并将所有未访问的子节点加入 `open_list`，从而实现 A* 算法。

2.4 实现 DQNRobot

Q-Learning 是一个值迭代 (Value Iteration) 算法。

与策略迭代 (Policy Iteration) 算法不同，值迭代算法会计算每个“状态”或是“状态-动作”的值 (Value) 或是效用 (Utility)，然后在执行动作的时候，会设法最大化这个值。因此，对每个状态值的准确估计，是值迭代算法的核心。通常会考虑最大化动作的长期奖励，即不仅考虑当前动作带来的奖励，还会考虑动作长远的奖励。

相较于传统根据 Q 值表估计 Q 值的方法，本报告采用 Deep-QLearning，以神经网络代替。用 PyTorch 实现：

```

1. class Robot(TorchRobot):
2.     def __init__(self, maze):
3.         super().__init__(maze)
4.         # 设迷宫奖励，负向强化目的地奖励以突出目标重要性
5.         maze.set_reward({"hit_wall": 10., "destination": -maze.maze_size ** 2 * 10, "default": 1.})
6.         self.maze = maze
7.         self.epsilon = 0
8.         # 建全图视野加速学习决策
9.         self.memory.build_full_view(maze)
10.        self.train()
11.
12.    def train(self):
13.        # 持续训练直至成功出迷宫
14.        while True:
15.            self._learn(batch=len(self.memory))
16.            self.reset()
17.            for _ in range(self.maze.maze_size ** 2):

```

```

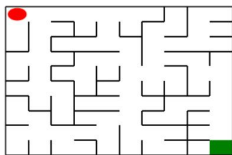
18.         action, reward = self.test
        _update()
19.         if reward == self.maze.reward["destination"]:
20.             return
21.
22.     def train_update(self):
23.         # 依状态选动作获奖励
24.         state = self.sense_state()
25.         action = self._choose_action(state)
26.         reward = self.maze.move_robot(action)
27.         return action, reward
28.
29.     def test_update(self):
30.         # 转状态为张量算 Q 值选动作获奖励
31.         state = np.array(self.sense_state(), dtype=np.int16)
32.         state = torch.from_numpy(state).float().to(self.device)
33.         with torch.no_grad():
34.             q_value = self.eval_model(state).cpu().data.numpy()
35.             action = self.valid_action[np.argmax(q_value).item()]
36.             reward = self.maze.move_robot(action)
37.             return action, reward

```

3. 结果分析

3.1 基础算法

在自己测试时，本报告采用了经典的 10*10 迷宫，效果良好。

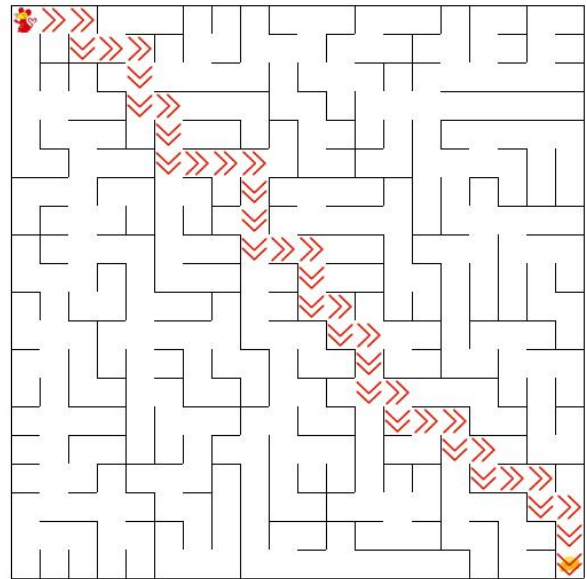


Maze of size (10, 10)
搜索出的路径: ['r', 'r', 'r', 'r', 'd', 'd', 'r', 'd', 'd', 'r', 'u', 'r', 'r', 'd', 'r', 'd', 'd', 'd', 'd', 'd']
恭喜你，到达了目标点

在实机测试时，其表现也非常优异，使用时长达到了惊人的 0s。

测试基础搜索算法	✓	0s	恭喜, 完成了迷宫
----------	---	----	-----------

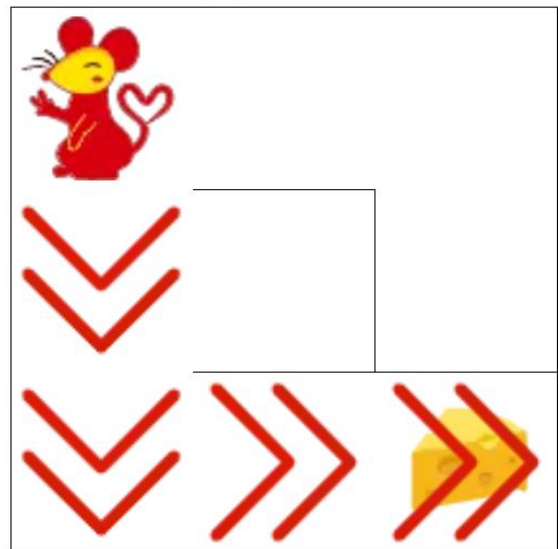
基础搜索算法 (Victory)



3.2 DQN

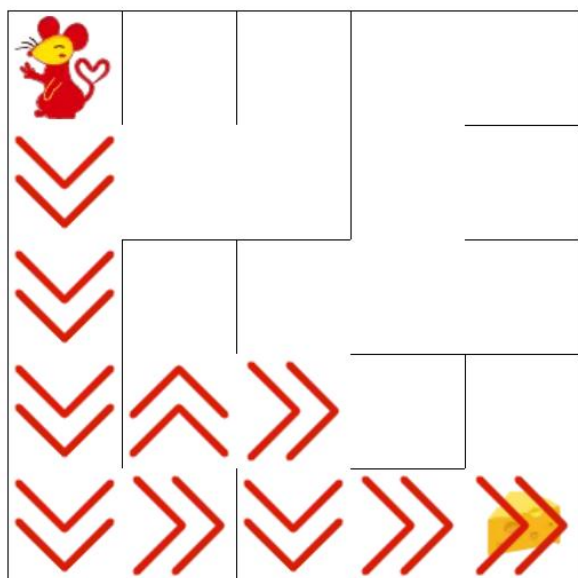
DQN 由于需要较长时间的训练过程，相对以上算法要慢很多。但是对于比较简单的迷宫，其仍然能够胜任：

强化学习level3 (Victory)



测试强化学习算法(初级)	✓	0s	恭喜, 完成了迷宫
--------------	---	----	-----------

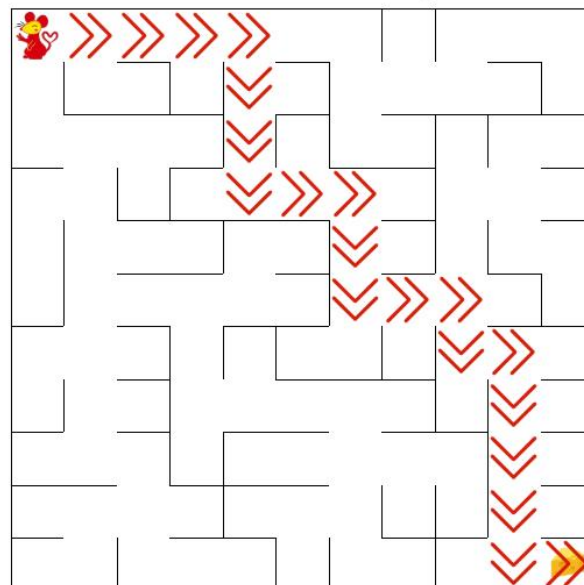
强化学习level5 (Victory)



测试强化学习 算法(中级)	✓	1s	恭喜, 完成了迷宫
------------------	---	----	-----------

而对于比较复杂的迷宫，其运行时间明显增大，达到了比较惊人的程度，不过最后还是完成了任务：

强化学习level11 (Victory)



测试强化学习 算法(高级)	✓	273s	恭喜, 完成了迷宫
------------------	---	------	-----------

4. 总结与讨论

针对 DQN 的算法问题，有数种改进方向。传统的 DQN 在评估 Q 值时可能会高估，从而导致不稳定的训练结果。使用 Double DQN 可以有效解决这个问题。方法为增加目标网络 (target_model)，使用目标网络来计算目标 Q 值，避免直接使用评估网络。也可以使用分布式强化学习，让多个代理并行探索迷宫并共享经验，可以大大加速训练。根据迷宫的复杂程度和运行环境，可以逐步实现这些改进方案，优先从最能提升性能的地方（如分布式 DQN）入手。