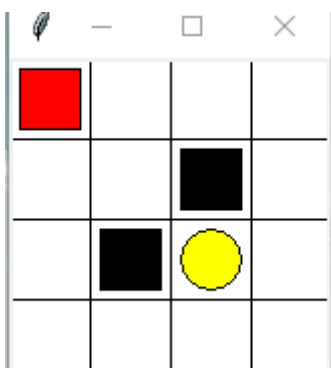


# 强化学习基础篇（十一） 迷宫环境搭建

这节中我们看看如何自己搭建一个强化学习实验环境，这里要做的是一个简单的迷宫环境。智能体在训练过程中的可视化过程如下：



## 1、环境设定

- 红色正方形表示在探索中的智能体
- 黑色正方形表示一个地狱终结点，当红色智能体到这个状态时，获得-1的奖励。
- 黄色位置表示天堂终结点，当红色智能体到这个状态时，获得+1的奖励。
- 所有其他的白色位置的奖励都为0

## 2. 源码信息如下

这里我们主要使用numpy, sys以及Tkinter。

```
1  # 导入库信息
2  import numpy as np
3  import time
4  import sys
5  if sys.version_info.major == 2:
6      import Tkinter as tk
7  else:
8      import tkinter as tk
9
10 # 设定环境信息
11 UNIT = 40 # 设定是像素大小为40
12 MAZE_H = 4 # 设置纵轴的格子数量
13 MAZE_W = 4 # 设置横轴的格子数量
14
15 # 创建一个迷宫类
16 class Maze(tk.Tk, object):
17     def __init__(self):
18         super(Maze, self).__init__()
19         # 定义动作空间为上下左右四个动作
20         self.action_space = ['u', 'd', 'l', 'r']
21         # 获取动作数量
22         self.n_actions = len(self.action_space)
23         # 定义迷宫名字
24         self.title('maze')
25         # 通过geometry函数来设置窗口的宽和高，分别为格子数量乘以像素大小
26         self.geometry('{0}x{1}'.format(MAZE_H * UNIT, MAZE_H * UNIT))
27         # 调用迷宫创建函数
```

```

28         self._build_maze()
29
30     def _build_maze(self):
31         # 设定画布大小
32         self.canvas = tk.Canvas(self, bg='white',
33                                 height=MAZE_H * UNIT,
34                                 width=MAZE_W * UNIT)
35
36         # 创建一个个的小格子
37         for c in range(0, MAZE_W * UNIT, UNIT):
38             x0, y0, x1, y1 = c, 0, c, MAZE_H * UNIT
39             self.canvas.create_line(x0, y0, x1, y1)
40         for r in range(0, MAZE_H * UNIT, UNIT):
41             x0, y0, x1, y1 = 0, r, MAZE_W * UNIT, r
42             self.canvas.create_line(x0, y0, x1, y1)
43
44         # 创建一个原点
45         origin = np.array([20, 20])
46
47         # 创建第一个地狱节点
48         hell1_center = origin + np.array([UNIT * 2, UNIT])
49         self.hell1 = self.canvas.create_rectangle(
50             hell1_center[0] - 15, hell1_center[1] - 15,
51             hell1_center[0] + 15, hell1_center[1] + 15,
52             fill='black')
53         # 创建第二个地狱节点
54         hell2_center = origin + np.array([UNIT, UNIT * 2])
55         self.hell2 = self.canvas.create_rectangle(
56             hell2_center[0] - 15, hell2_center[1] - 15,
57             hell2_center[0] + 15, hell2_center[1] + 15,
58             fill='black')
59
60         # 创建一个圆形的天堂节点
61         oval_center = origin + UNIT * 2
62         self.oval = self.canvas.create_oval(
63             oval_center[0] - 15, oval_center[1] - 15,
64             oval_center[0] + 15, oval_center[1] + 15,
65             fill='yellow')
66
67         # 创建红色探索得智能体
68         self.rect = self.canvas.create_rectangle(
69             origin[0] - 15, origin[1] - 15,
70             origin[0] + 15, origin[1] + 15,
71             fill='red')
72
73         # ppack函数的作用是让画布显示中正确的位置上。如果没调用这个函数，就不会正常地显示任何东西。
74         self.canvas.pack()
75         # 定义重置函数
76         def reset(self):
77             self.update()
78             # 设定重置函数延迟0.5秒
79             time.sleep(0.5)
80             # 重置的时候删除原来智能体位置，然后重新设定其在原点
81             self.canvas.delete(self.rect)
82             origin = np.array([20, 20])
83             self.rect = self.canvas.create_rectangle(
84                 origin[0] - 15, origin[1] - 15,

```

```

85         origin[0] + 15, origin[1] + 15,
86         fill='red')
87     # 返回重置后的位置
88     return self.canvas.coords(self.rect)
89 # 定义每步探索函数
90 def step(self, action):
91     # 首先获取当前的坐标
92     s = self.canvas.coords(self.rect)
93     # 根据动作来定义智能体会如何移动
94     base_action = np.array([0, 0])
95     if action == 0: # 向上移动的情况
96         if s[1] > UNIT:
97             base_action[1] -= UNIT
98     elif action == 1: # 向下移动的情况
99         if s[1] < (MAZE_H - 1) * UNIT:
100             base_action[1] += UNIT
101     elif action == 2: # 向右移动的情况
102         if s[0] < (MAZE_W - 1) * UNIT:
103             base_action[0] += UNIT
104     elif action == 3: # 向左移动的情况
105         if s[0] > UNIT:
106             base_action[0] -= UNIT
107     # 按照动作移动智能体
108     self.canvas.move(self.rect, base_action[0], base_action[1])
109     # 移动后的状态定义为s_
110     s_ = self.canvas.coords(self.rect)
111
112     # 根据移动后的下一个状态，设定奖励值。
113     if s_ == self.canvas.coords(self.oval):
114         reward = 1
115         done = True
116         s_ = 'terminal'
117     elif s_ in [self.canvas.coords(self.hell1),
118 self.canvas.coords(self.hell2)]:
119         reward = -1
120         done = True
121         s_ = 'terminal'
122     else:
123         reward = 0
124         done = False
125         info = None
126     # 返回下一个状态，奖励以及当前Episode是否完成的
127     return s_, reward, done, info
128 # 刷新当前环境
129 def render(self):
130     time.sleep(0.1)
131     self.update()
132
133 # 刷新函数
134 def update():
135     for t in range(10):
136         s = env.reset()
137         while True:
138             env.render()
139             a = 1
140             s, r, done, info = env.step(a)
141             if done:
142                 break

```

```

142 # 自测部分
143 if __name__ == '__main__':
144     env = Maze()
145     env.after(100, update)
146     env.mainloop()

```

### 3、说明

这个迷宫环境的创建虽然很简单，但是还是按照Gym的框架来搭建的。他包含了一个强化学习环境应有的几个要素，包括以下主要几个函数：

- reset(self):重置环境的状态，返回环境的当前状态。
- step(self, action): 推进一个时间步，返回信息是state, reward, done, info。其中，done用来表明游戏是否结束的标志位，Info是关于游戏的附加信息。
- render (self) : 重新绘制环境

### 4、环境的调用

这里算法的核心后续再介绍，主要看看在运行中是如何对环境进行调用的。

```

1  def train():
2      for episode in range(100):
3          # 初始化观测
4          observation = env.reset()
5
6          while True:
7              # 刷新环境
8              env.render()
9
10             # 先通过强化学习算法按照既定策略选择一个动作
11             action = RL.choose_action(str(observation))
12
13             # 根据选择的动作，调用环境返回下一个状态，奖励以及结束标志位
14             observation_, reward, done = env.step(action)
15
16             # 调用强化学习算法的训练过程
17             RL.learn(str(observation), action, reward, str(observation_))
18
19             # 将当前获取的下一步状态设定为新的状态，并准备下一个循环
20             observation = observation_
21
22             # 按照标志位结束while
23             if done:
24                 break
25
26             # 游戏结束
27             print('game over')
28             env.destroy()
29
30 if __name__ == "__main__":
31     env = Maze()
32     RL = QLearningTable(actions=list(range(env.n_actions)))
33
34     env.after(100, train)
35     env.mainloop()

```

环境相关的信息就介绍到这里，后续该开始算法部分了。

---

## 历史文章链接：

- [强化学习基础篇（十） OpenAI Gym环境汇总](#)
- [强化学习基础篇（九） OpenAI Gym基础介绍](#)
- [强化学习基础篇（八） 动态规划扩展](#)
- [强化学习基础篇（七） 动态规划之价值迭代](#)
- [强化学习基础篇（六） 动态规划之策略迭代（2）](#)
- [强化学习基础篇（五） 动态规划之策略迭代（1）](#)
- [强化学习基础篇（四） 动态规划之迭代策略评估](#)
- [强化学习基础篇（三） 动态规划之基础介绍](#)
- [强化学习基础篇（二） 马尔科夫决策过程（MDP）](#)
- [强化学习基础篇（一） 强化学习入门](#)
- [9.进一步讨论Policy Gradients方法](#)
- [8. DRL中的Q-Function](#)
- [7. 值函数方法（Value Function Methods）](#)
- [6. Actor-Critic算法](#)
- [5. 策略梯度（Policy Gradients）](#)
- [4.强化学习简介](#)
- [3.TensorFlow示例](#)
- [2.模仿学习（Imitation Learning）](#)
- [1.深度强化学习简介](#)
- [A survey on value-based deep reinforcement learning](#)
- [Chinese Stock Prediction Using Deep Neural Network](#)
- [Differential Dynamics of the Maternal Immune System](#)阅读笔记
- [bib如何生成author-year格式的bbl的问题](#)
- [如何使用GPU运行TensorFlow（Win10）](#)
- [通过frp实现内网穿透](#)
- [关于t-SNE降维方法](#)
- [系统评价与Meta分析基础](#)
- [Meta分析入门工具介绍](#)
- [使用ruptures检测变量关系](#)
- [如何读取rda格式数据](#)