

Tic Tac Toe

于11月 26, 2018由SWRhapsody发布

Introduction

刚开学这一个月真是好忙啊，我都没什么时间来写博客。这几天看了点 Reinforcement Learning 相关的知识写个小总结。因为我自己的专业不是这方面所以为了避免瞎讲，我无法在这篇博客涉及太多的概念。大部分概念可以在 [Reference](#) 中看。

Code

Tic Tac Toe!

Reinforcement Learning是机器学习的一部分，这篇博文的目标是用 Reinforcement Learning 实现个玩 Tic Tac Toe 的游戏AI，完整的代码[\[7\]](#)，我大部分的代码都是参照[\[2\]](#)。

我们先写个 Board 来作为一个游戏场地

```
1 import numpy as np
2 import pandas as pd
3
4 class Game:
5     def __init__(self, row, col):
6         self.row_size = row
7         self.col_size = col
8         # 0 for not end
9         # 1 for winner is first player
10        # 2 for winner is second player
11        # 3 for tie
12        self.end = 0
13        self.board = np.zeros((self.row_size, self.col_size))
14        self.size = self.board.size
15        self.mapping = {0: " ", 1: "X", -1: "O"}
16
17
18    def DrawCharForItem(self, item):
19        # 0 for ' '
20        # 1 for 'X'
21        # -1 for 'O'
22        return self.mapping[item]
```

SWRhapsody

```

27     board =
28         {} | {} | {}
29         -----
30         {} | {} | {}
31         -----
32         {} | {} | {}
33     """.format(*items)
34     print(board)
35
36     def CheckWin(self):
37         result = []
38
39         for i in range(self.row_size):
40             result.append(np.sum(self.board[i, :]))
41         for i in range(self.col_size):
42             result.append(np.sum(self.board[:, i]))
43
44         result.append(0)
45         for i in range(0, self.col_size):
46             result[-1] += self.board.item((i, i))
47         result.append(0)
48         for i in range(0, self.row_size):
49             result[-1] += self.board.item((i, self.row_size - 1 - i))
50
51         for i in result:
52             if i == 3:
53                 self.end = 1
54                 return self.end
55             if i == -3:
56                 self.end = 2
57                 return self.end
58         sum = np.sum(np.abs(self.board))
59         if sum == self.size:
60             self.end = 3
61         return self.end
62
63     def End(self):
64         return self.end
65
66     def Step(self, order, action):
67         if action[0] < 0 or action[0] > self.row_size-1 or action[1] < 0 or action[1] > self.col_size-1:
68             print("x or y invalid")
69             return True
70         if self.board[action[0], action[1]] != 0:
71             print("this palce has already been taken")
72             return True
73         self.board[action[0], action[1]] = order
74         self.CheckWin()
75         self.Print()
76         return False

```

接下来写个Agent类作为我们的AI。

```
1 class Agent
```

现在来介绍下我们这个AI的核心公式 ([8])



SWRhapsody

不过在编码这公式之前我们需要一个数据结构来存储我们的收益。

定义一个 `state` 来存储收益。这个 `state` 是一个 `row*col` 的数组，这个数组存放当前情况下每一格的收益。

```
1 state, 当前情况下走 (0,1) 这个点的收益是 1.30938546e-07
2 [[ 1.30938546e-07 -4.93034771e-07 1.51138875e-04]
3  [-6.41862293e-08 -2.77494609e-08 -5.18934185e-07]
4  [-2.61877061e-07 -4.96284826e-05 2.59479593e-07]]
```

我们在定义一个字典来存每一种情况下的 `state`，为每一种情况生成一个 `hash`，这个 `hash` 作为字典的 `key`

```
1 def ComputeHash(board):
2     board_ = board.flatten().tolist()
3     return "".join([str(i) for i in board_])
4
5 class Agent:
6     def __init__(self, name, exploration_rate=0.33, learning_rate=0.5, discount_factor=0.01):
7         self.states = {}
8         # states stack
9         self.state_order = []
10        self.learning_rate = learning_rate
11        self.exploration_rate = exploration_rate
12        self.discount_factor = discount_factor
13        self.name = name
```

这样我们的公式就变成了这样

```
1 class Agent:
2     def learn_by_temporal_difference(self, reward, new_state_hash, state_hash):
3         prev_state = self.states.get(
4             state_hash, np.zeros((BOARD_ROWS, BOARD_COLS)))
5         return self.learning_rate*((reward*self.states[new_state_hash] )- prev_state)
```

接着定义一个 `reward` 函数，对于什么时候给予奖励有多种选择[9]，在这个 AI 中我们选择在出现赢家的時候给予奖励。

```
1 class Agent:
2     def set_state(self, board, action):
3         hash = ComputeHash(board)
4         self.state_order.append((hash, action))
5
6     def OnReward(self, reward):
7         hash, action = self.state_order.pop()
8         if hash not in self.states:
9             self.states[hash] = np.zeros((BOARD_ROWS, BOARD_COLS))
10        self.states[hash].itemset(action, reward)
11
12        while self.state_order:
13            hash_prev, action_prev = self.state_order.pop()
14            reward *= self.discount_factor
```

SWRhapsody

```

20         self.states[hash_prev] = np.zeros((BOARD_ROWS, BOARD_COLS))
21         reward = self.learn_by_temporal_difference(reward, hash, hash_prev).item(action)
22         self.states[hash_prev].itemset(action, reward)
23
24         hash = hash_prev
25         action = action_prev

```

这里 `state_order` 是我们存所走过的每一步的一个栈，存的是 (state的hash, 这一步的坐标)。这个函数的组要目的是把计算出的收益存储在 `self.states` 中。

接下来我们定义 `exploit` 和 `explore` 函数

```

1 class Agent:
2     def exploit(self, board):
3         state_value = self.states[ComputeHash(board)]
4         x, y = np.where(state_value == state_value.max())
5         best_choices = [(a, b) for a, b in zip(x, y)]
6         return best_choices[np.random.choice(len(best_choices))]
7
8     def explore(self, board):
9         x, y = np.where(board == 0)
10        vacant = [(a, b) for a, b in zip(x, y)]
11        return vacant[np.random.choice(len(vacant))]

```

最后定义个用来选择走那一步的函数，何时选择探索具体可以参考[5]中的2.2

```

1 class Agent:
2     def SelectMove(self, board):
3         action = None
4         exploration = np.random.random() < self.exploration_rate
5         if exploration or hash not in self.states:
6             print("%s exploit" % self.name)
7             action = self.explore(board)
8         else:
9             print("%s exploit" % self.name)
10            action = self.exploit(board)
11            # update state
12            self.set_state(board, action)
13            return action

```

训练下这个AI

```

1 def Train(round, bot1, bot2):
2     win_trace = pd.DataFrame(data=np.zeros(
3         (round, 2)), columns=["bot1", "bot2"])
4     for i in range(round):
5         print("-"*100)
6         print("Round:{}".format(i+1))
7         game = Game(BOARD_ROWS, BOARD_COLS)
8         turn = 1
9         while game.End() == 0:
10            if turn == 1:
11                action = bot1.SelectMove(game.board)
12                game.Step(1, action)
13                turn = 2
14            else:
15                action = bot2.SelectMove(game.board)

```

SWRhapsody

```

21         bot2.OnReward(-1)
22         win_trace.set_value(i, 'bot1', 1)
23     elif game.End() == 2:
24         bot1.OnReward(-1)
25         bot2.OnReward(1)
26         win_trace.set_value(i, 'bot2', 1)
27     return win_trace
28
29 if __name__ == "__main__":
30     bot1 = Agent("bot1")
31     bot2 = Agent("bot2")
32     Train(5000, bot1, bot2)

```

最后跟它玩下

```

1 Bot first!
2 bot1 exploit
3
4         |   | X
5         -----
6         |   |
7         -----
8         |   |
9
10 Your turn(enter x,y):1,1
11
12         |   | X
13         -----
14         | 0 |
15         -----
16         |   |
17
18 bot1 exploit
19
20         |   | X
21         -----
22         | 0 |
23         -----
24         |   | X
25
26 Your turn(enter x,y):1,2
27
28         |   | X
29         -----
30         | 0 | 0
31         -----
32         |   | X
33
34 bot1 exploit
35
36         |   | X
37         -----
38         X | 0 | 0
39         -----
40         |   | X
41
42 Your turn(enter x,y):0,1
43
44         | 0 | X
45         -----

```

SWRhapsody

```

51
52         | 0 | X
53         -----
54        X | 0 | 0
55         -----
56        X |   | X
57
58 Your turn(enter x,y):2,1
59
60         | 0 | X
61         -----
62        X | 0 | 0
63         -----
64        X | 0 | X
65
66 You win

```

尝试了比较多的次数，只有一次训练出来的智商过关(结果忘保存结果)。感觉训练的次数不够并且需要把规则教给AI，不然训练的速度太慢。

Reference

- [1] https://en.wikipedia.org/wiki/Reinforcement_learning
- [2] https://github.com/AmreshVenugopal/tic_tac_toe/blob/master/Tic%20tac%20toe.ipynb
- [3] https://github.com/ShangtongZhang/reinforcement-learning-an-introduction/blob/master/chapter01/tic_tac_toe.py
- [4] https://en.wikipedia.org/wiki/Temporal_difference_learning
- [5] <http://tokic.com/www/tokicm/publikationen/papers/AdaptiveEpsilonGreedyExploration.pdf>
- [6] <https://detailed.af/a-game-of-tic-tac-toe/>
- [7] https://github.com/lceware/blog_code/blob/master/2018/tic_tac_toe.py
- [8] <http://incompleteideas.net/book/the-book.html>
- [9] <https://medium.freecodecamp.org/an-introduction-to-reinforcement-learning-4339519de419>

分类: **UNCATEGORIZED**



SWRhapsody

0 条评论

发表评论

名称 *

电子邮件 *

网站

在想些什么?

发表评论

近期文章

[携程Apollo YAML 反序列化](#)[CVE-2020-5410](#)[CodeQL部分源码简读](#)[服务器与网关的不一致](#)

SWRhapsody

近期评论

文章归档

2020年8月

2020年6月

2020年5月

2020年3月

2020年1月

2019年12月

2019年11月

2019年8月

2019年7月

2019年5月

2019年4月

2019年1月

2018年11月

2018年10月

2018年9月

2018年4月

2018年3月

2018年2月

2018年1月

分类目录

SWRhapsody

[cryptography](#)[Exercise](#)[Exploit](#)[HackTheBox](#)[Penetration Test](#)[Uncategorized](#)

相关文章

UNCATEGORIZED

WordPress Social Warfare XSS

Introduction Plugin name: Social Warfare [阅读更多...](#)

UNCATEGORIZED

自动化程序修复

Automatic patch generation Introduction [阅读更多...](#)

SWRhapsody

Criminology

Introduction 网络安全也要学习犯罪学了吗，每次写这类文章都 [阅读更多...](#)

ABOUT

Hestia | 由Themelsle开发