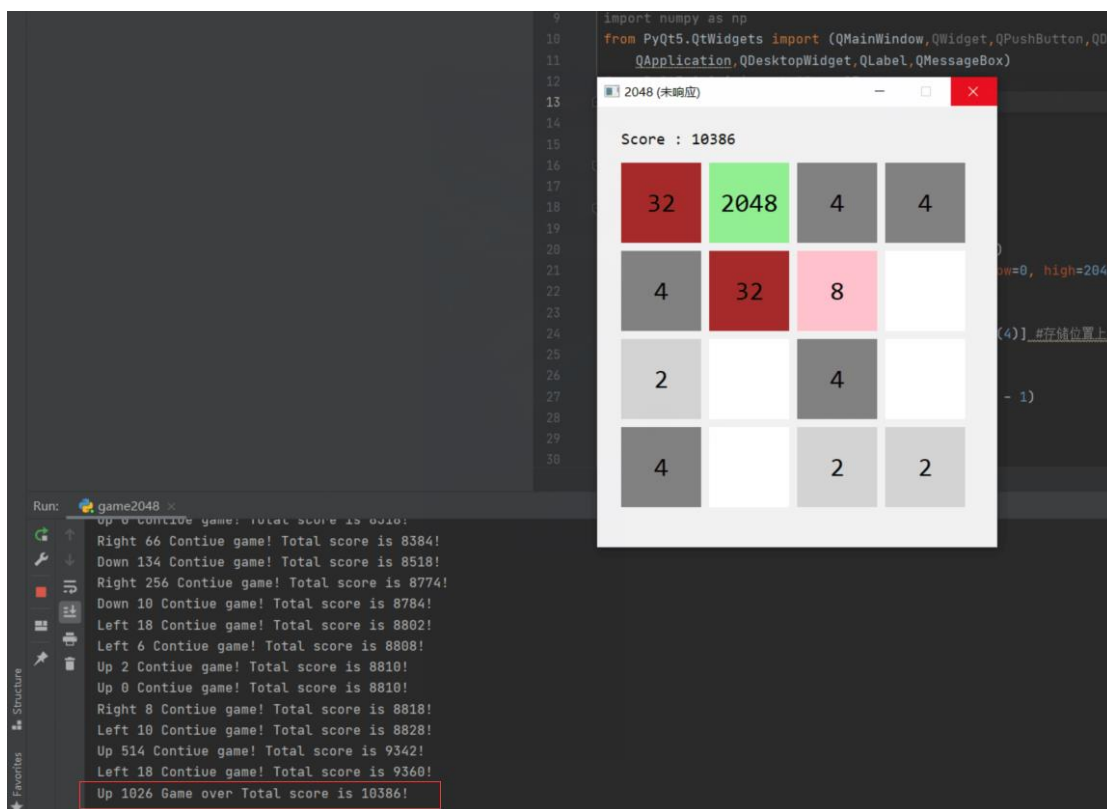


PA2 实验报告

191300073 杨斯凡 191300073@smail.nju.edu.cn

一、实验结果



我注释掉了结束后关闭的代码，就可以不退出。

二、具体思路

这次的问题是 2048 游戏，框架代码将游戏限制到 4x4 的方格上，玩家可以采取上下左右四个行动，相同的数字格子会合并。但是，当采取行动之后，下一个 2 的生成位置是随机的，这就导致采取动作之后的下一个状态未知。

而在阅读完框架代码后，我对给定的接口 API 有了一个初步的了解，之后我使用了课上讲的蒙特卡洛树搜索来进行最佳行动的选取。

我定义了两个类，一个是蒙特卡洛树的节点类，存储了每个节点的状态，Q，N，这个节点的子节点和父节点。另一个是整个蒙特卡洛树，其中各个函数的作用如下：

calculate_value: 计算每个节点的上置信界，如果没有探索，则会返回正无穷（使用了 $1e50$ 模拟正无穷）。

expand_choose: 是算法中的展开，选择两步，首先选择行动，然后把对应行动的子节点展开。

Evaluate: 使用了 rollout 方法，默认策略是随机策略，到达预设深度（10）后使用当前评分进行评价。

Reverse_update: 沿着当前节点的父节点向上直到根节点，沿路更新节点的 Q 和 N。

三、具体实现

计算上置信界，使用了 PPT 上的公式，如果未探索则返回正无穷，其中我将 c 设置为 100。

```

358     def caculate_value(self,node,action):
359         if node.child[action] == None:
360             return infinity
361         else:
362             return node.child[action].Q + c * math.sqrt(math.log(node.N) / node.child[action].N)

```

$$Q(s,a) + c \sqrt{\frac{\log N(s)}{N(s,a)}}$$

评价：我使用了均匀随机的 rollout 方法，深度设置为 10，达到深度 10 之后，选择当前的得分作为值返回。

```

385     def evalute(self,node,env):
386         value=env.score
387         for i in range(d):
388             action = random.randint(0, 3)
389             obs, rew, done, info = env.step(action)
390             if done==True:
391                 break
392             value = env.score
393         return value

```

反向更新：从当前节点一直向上遍历，直到遍历到根节点位置，更新每个节点的 Q 和 N

```

394     def Reverse_update(self,node,value):
395         while node != None:
396             node.N += 1
397             node.Q = node.Q + (value - node.Q) / node.N
398             node = node.father
399
400

```

展开选择：首先拷贝当前的所有状态进行计算，然后从根节点开始进行选择最佳的动作，每次模拟做出行动后的回报，直到叶节点返回当前状态。

```

363     def expand_choose(self):#因为初始状态已经初始化过,所以先选择
364         temp=copy.deepcopy(env)
365         node=self.root
366
367         while 1:
368             action=0
369             max_value = -infinity
370             for i in range(4):
371                 now_value=self.caculate_value(node,i)
372                 if now_value > max_value:
373                     action=i
374                     max_value=now_value
375
376             if node.child[action]==None:
377                 obs, rew, done, info = temp.step(action)
378                 child_node = Node(temp.state)
379                 child_node.father = node
380                 node.child[action] = child_node
381                 return child_node, temp
382             else:
383                 node = node.child[action]
384                 obs, rew, done, info = temp.step(action)

```

主循环:

首先获取当前的时间, 以保证 1s 内做出选择, 每次搜索重置整个搜索树, 以当前状态作为新的根节点搜索。循环体在 1s 内会完成选择, 展开, 评估和反向更新。然后再从根节点的子节点中选择值最大的行动以作为最优行动。

```

401     def Play(self,state,env):
402         self.root = Node(state)
403         self.env = copy.deepcopy(env)
404         t = time.time()
405         t=int(round(t * 1000))
406         search_time=0
407         while time_bound > int(round(time.time() * 1000))-t:
408             new_node, new_env = self.expand_choose()
409             value = self.evalute(new_node, new_env)
410             self.Reverse_update(new_node, value)
411             action = 0
412             max_value = -infinity
413
414             for i in range(4):
415                 if self.root.child[i] != None:
416                     value = self.root.child[i].Q
417                     if value > max_value:
418                         max_value = value
419                         action = i
420
421             self.env.step(action)
422             self.root = self.root.child[action]
423             self.root.father = None
424             return action
425

```

但是在实验的过程中有很多的 bug。

首先是在快产生 2048 的时候，会报错 `KeyError`，这个是框架代码的问题，和助教沟通后，我将代码进行了修改

```
col = [lightgrey, grey, pink, yellow, brown,
for num in range(1, 15):
    self.color_dict[2 ** num] = col[num - 1]
```

相应的颜色也做出了修改。

第二个问题就是，右上角的得分只会显示 2 位，和助教沟通后发现是 linux 和 windows 的问题，我将字体调小后，就可以显示 4 位。

四、运行方式

直接运行即可