

# 强化学习 HW1

— 人 饶翔云 520030910366

## Problem 1:

首先以数学形式重新阐述问题：

1. If we produce a policy  $\pi'$  from  $\pi$ , then  $V^{\pi'} \geq V^\pi$  (因为如果该公式成立, 那么由于策略 $\pi'$ 是根据 $V^\pi$ 得来的, 而 $V^{\pi'} \geq V^\pi$ , 由 $V^{\pi'}$ 所得到的策略肯定不比 $\pi'$ 差)
2. If the optimal policy is  $\pi_{opt}$ , our policy will converges to it finally, and that means  $V^{\pi_{opt}}$  is optimal with  $V^\pi$  will finally converges to it.(最优策略对应最大价值函数)

对于第一个问题：

当价值函数收敛后, 对于任意状态 $s$ :

$$\begin{aligned} V^{\pi'}(s) &= \sum_{s'} P(s, \pi'(s), s') [R(s, \pi'(s), s') + \gamma V^{\pi'}(s')] = R(s, \pi'(s)) + \gamma \sum_{s'} P(s, \pi'(s), s') V^{\pi'}(s') \\ V^\pi(s) &= \sum_{s'} P(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')] = R(s, \pi(s)) + \gamma \sum_{s'} P(s, \pi(s), s') V^\pi(s') \end{aligned}$$

从迭代公式易知:  $V_0^{\pi'} = V^\pi$

又因为 $\pi'(s) = \operatorname{argmax}_a \sum_{s'} P(s, a, s') [R(s, a, s') + \gamma V^\pi(s')]$ , 所以:

$$V^\pi(s) \leq R(s, \pi'(s)) + \gamma \sum_{s'} P(s, \pi'(s), s') V_0^{\pi'}(s') = V_1^{\pi'}(s)$$

故 $V_1^{\pi'}(s) - V_0^{\pi'}(s) \geq 0$ . 假设已知 $V_i^{\pi'}(s) - V_{i-1}^{\pi'}(s) \geq 0, \forall s$ , 那么:

$$V_{i+1}^{\pi'}(s) - V_i^{\pi'}(s) = \gamma \sum_{s'} P(s, \pi', s') (V_i^{\pi'}(s') - V_{i-1}^{\pi'}(s')) \geq 0, \forall s$$

所以 $V_i^{\pi'}(s) \leq V_{i+1}^{\pi'}(s), \forall s \in S, \forall i \in \mathbb{N}$ ,

因此 $V^{\pi'}(s) = V_\infty^{\pi'}(s) \geq V_0^{\pi'}(s) = V^\pi(s)$ , 第一个问题证毕。

对于第二个问题：

$$\forall i \in \mathbb{N}, \forall s \in S, |V_{\pi_{i+1}} - V_{\pi_i}| \leq \gamma |V_{\pi_i} - V_{\pi_{i-1}}| \text{待证}$$

奖励函数有界, 不妨令 $R(s, a) \leq n$

$$\text{那么 } V^\pi = \mathbb{E}(\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s'_t) | \pi, s_0 = s) \leq n(1 + \gamma + \gamma^2 + \dots) = \frac{n}{1-\gamma}$$

将 $\{V^{\pi_i}\}$ 视为一个数列, 那么根据单调递增有界收敛定理,  $\{V^{\pi_i}\}$ 必定收敛。

那么 $\{V^{\pi_i}\}$ 必有上确界, 而这就是我们所求的最大价值函数, 由这个最大价值函数得到的策略就是最优策略。

## problem 2:

---

首先列举这个环境的性质:

- 二维的栅格世界
- 智能体有20%概率采取与理论不同的行为 ( $\epsilon$ -policy?)
- 智能体的目标是进入 `TERMINAL STATE` 并采取 `exit` 从而获得最终的奖励

### (a) 实现基于价值迭代的智能体

共有三个部分的代码需要补全。

#### 1. `runValueIteration()`

```
1  # implement the synchronous value iteration
2      for _ in range(0, self.iterations):
3          new_values = util.Counter()
4          for state in self.mdp.getStates():
5              # If the state is terminal, then the value is 0
6              if self.mdp.isTerminal(state):
7                  new_values[state] = 0
8              else:
9                  # Otherwise, the value is the maximum Q of all possible
10                 # actions(using old values)
11                 actions = self.mdp.getPossibleActions(state)
12                 new_values[state] = max(
13                     [self.computeQValueFromValues(state, action)
14                      for action in actions])
15             # If the maximum change is less than epsilon, then stop
16             change = [abs(self.values[state] - new_values[state])
17                       for state in self.mdp.getStates()]
18             if max(change) < self.epsilon:
19                 break
20             # Update the values
21             self.values = new_values
```

## 2. computeQValueFromValues()

```
1 def computeQValueFromValues(self, state, action):
2     """Compute the Q-value of action in state from the value function
3     stored in self.values."""
4
5     value = None
6
7     """ YOUR CODE HERE """
8     # The Q-value is the sum of all possible next states and
9     # their probabilities
10    value = 0
11    for next_state, prob in self.mdp.getTransitionStatesAndProbs(state, action):
12        value += prob * (self.mdp.getReward(state, action, next_state)
13                        + self.discount * self.values[next_state])
14
15    return value
```

## 3. computeActionFromValues()

```
1 def computeActionFromValues(self, state):
2     """The policy is the best action in the given state
3     according to the values currently stored in self.values.
4
5     You may break ties any way you see fit. Note that if
6     there are no legal actions, which is the case at the
7     terminal state, you should return None.
8     """
9
10    bestaction = None
11
12    """ YOUR CODE HERE """
13    # If the state is terminal, then there is no action
14    if not self.mdp.isTerminal(state):
15        # Find the action with the maximum Q
16        actions = self.mdp.getPossibleActions(state)
17        bestaction = max(actions, key=lambda action:
18                        self.computeQValueFromValues(state, action))
19
20    return bestaction
21
```

## (b) 实现基于策略迭代的智能体

共有三个部分的代码需要补全。

## 1. runValueIteration()

```
1 def runPolicyIteration(self):
2     """ YOUR CODE HERE """
3     # calculate V according to the policy
4     # initialize the policy
5     for state in self.mdp.getStates():
6         # If the state is terminal, then the policy is None
7         if self.mdp.isTerminal(state):
8             self.policy[state] = None
9         else:
10            # Otherwise, the policy is the one of the possible actions
11            actions = self.mdp.getPossibleActions(state)
12            idx = random.randint(0, len(actions)-1)
13            self.policy[state] = actions[idx]
14    current_policy = self.policy
15    for _ in range(0, self.iterations):
16        while True:
17            new_values = util.Counter()
18            for state in self.mdp.getStates():
19                # If the state is terminal, then the value is 0
20                if self.mdp.isTerminal(state):
21                    new_values[state] = 0
22                else:
23                    # Otherwise, the value is the maximum Q of all possible
24                    # actions(using old values)
25                    action = current_policy[state]
26                    new_values[state] =
27                        self.computeQValueFromValues(state, action)
28                # If the maximum change is less than epsilon, then stop
29                change = [abs(self.values[state] - new_values[state]) for state
30                        in self.mdp.getStates()]
31                # Update the values
32                self.values = new_values
33                if max(change) < self.epsilon:
34                    break
35            # calculate the new policy
36            new_policy = dict()
37            for state in self.mdp.getStates():
38                # If the state is terminal, then the policy is None
39                if self.mdp.isTerminal(state):
40                    new_policy[state] = None
41                else:
42                    # Otherwise, the policy is the action with the maximum Q
43                    actions = self.mdp.getPossibleActions(state)
44                    new_policy[state] = max(actions,
45                        key=lambda action: self.computeQValueFromValues(state, action))
46            # If the new policy is the same as the old policy, then stop
47            if new_policy == current_policy:
48                break
49            else:
50                current_policy = new_policy
51    # Update the policy
52    self.policy = current_policy
```

2. computeQValueFromValues()

同基于价值迭代的智能体。

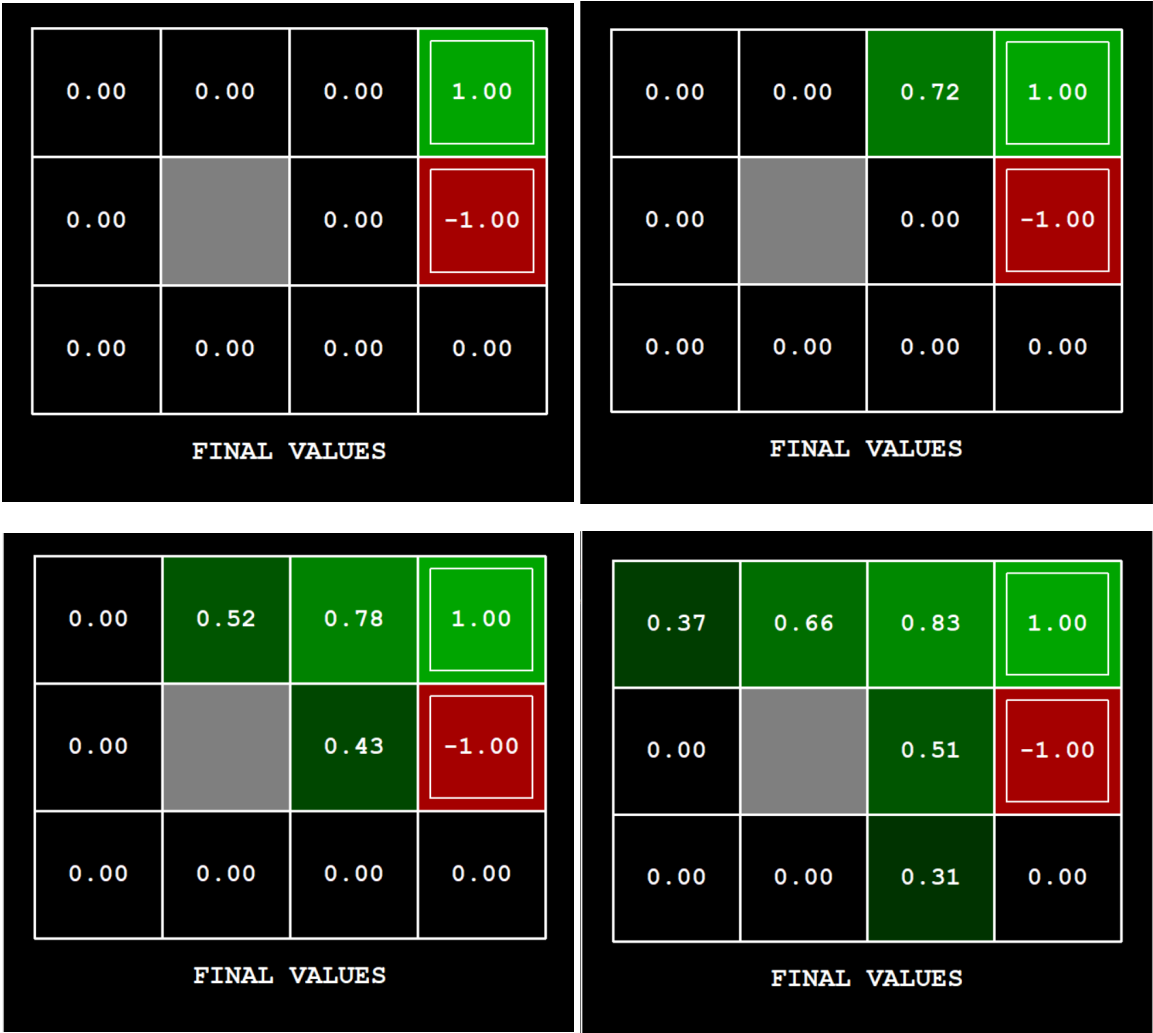
3. computeActionFromValues()

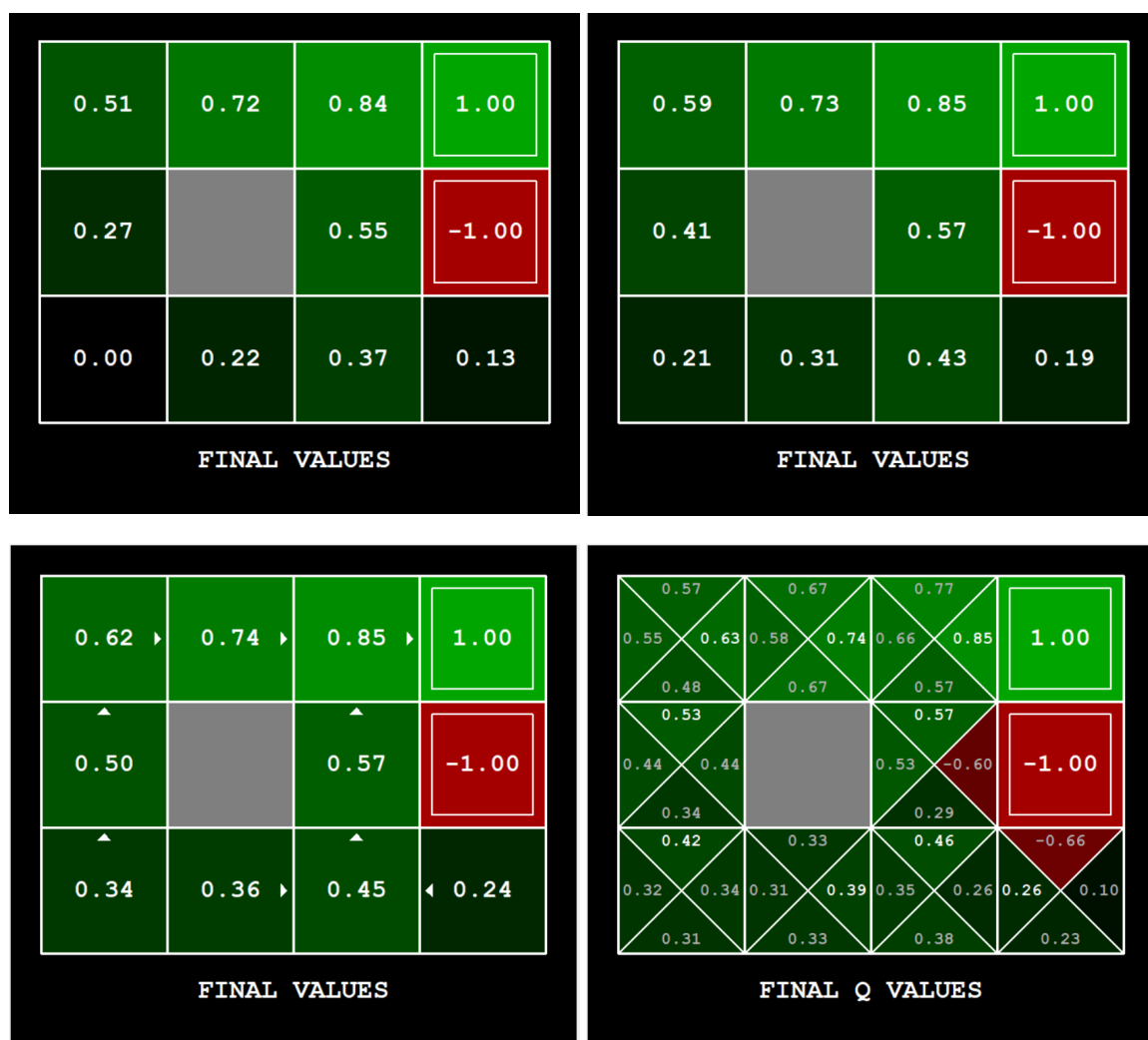
同基于价值迭代的智能体。

(c) 画图，探究收敛速度并说出原因

我将 $\epsilon$ 改为0.1，理由是更能看出收敛速度的差别。

价值迭代收敛后结果如下：





经过探究，V和Q在第7轮完全收敛，花了0.002秒（在0.001~0.003秒浮动）

策略迭代收敛后结果如下：

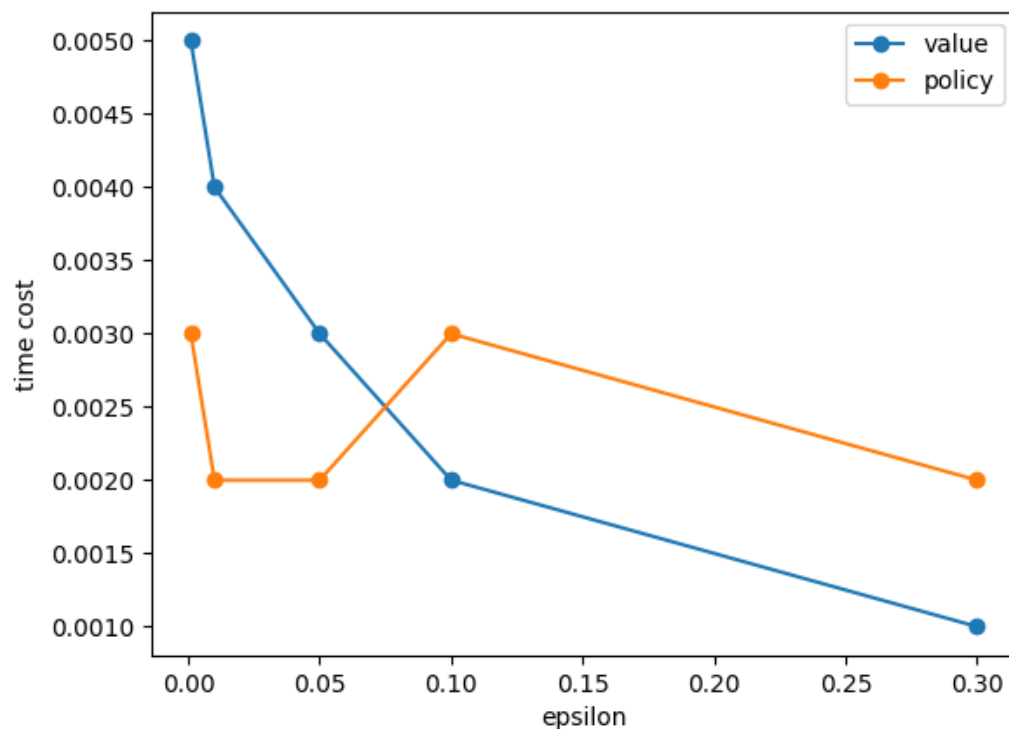




经过探究，V和Q在第4轮就已经完全收敛，花了0.003秒（在0.002~0.004秒浮动）

将 $\epsilon$ 改为0.01后，发现策略迭代需要0.002秒，而价值迭代需要0.004秒。

设置 $\epsilon \in [0.001, 0.01, 0.05, 0.1, 0.3]$ ,画出统计图如下:



**结论:** 在这个gridworld中, 当 $\epsilon$ 比较小的时候, 策略迭代收敛速度更快; 当 $\epsilon$ 比较大的时候, 价值迭代更快。

**原因:** 这个结论的前提是在这个gridworld中, 因为这个gridworld相对较小, 很容易找到最优策略, 策略迭代和价值迭代时间开销才可比较。因为策略迭代在每一轮迭代都会进行一次价值迭代 (基于当前策略)。当 $\epsilon$ 比较小的时候, 价值函数要收敛花费时间长, 价值迭代轮数多。但是策略迭代的每一轮都将价值迭代中贪心求最优动作的步骤改为基于当前策略取动作, 这一步节省的开销由于迭代轮数变多从而减小了总时间开销。当 $\epsilon$ 比较大的时候, 价值迭代轮数少, 策略迭代每轮的价值迭代所节约的开销就减少了, 而要策略要迭代的轮数还在那, 从而时间开销比正常的价值迭代更大。