

强化学习基础篇（十四）价值迭代在FrozenLake中的实现

本节将主要基于gym环境中的FrozenLake-v0进行价值迭代的实现。

1. 价值迭代算法的伪代码

价值迭代算法，用于估计 $\pi = \pi_*$

算法参数：小阈值 $\theta > 0$ ，用于确定估计量的精度。

对于任意 $s \in S^+$ ，任意初始化 $V(s)$ ，其中 $V(\text{终止状态}) = 0$

循环：

$\Delta \leftarrow 0$

对每一个 $s \in S$ 循环：

$v \leftarrow V(s)$

$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

直到 $\Delta < \theta$

输出一个确定的 $\pi \approx \pi_*$ ，使得

$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

2. 源代码

```
1 import numpy as np
2
3 def calc_action_value(state, v, discount_factor=1.0):
4     """
5     Calculate the expected value of each action in a given state.
6     对于给定的状态 s 计算其动作 a 的期望值
7     """
8     A = np.zeros(env.nA)
9     for a in range(env.nA):
10         for prob, next_state, reward, done in env.P[state][a]:
11             A[a] += prob * (reward + discount_factor * v[next_state])
12     return A
13
14 def value_iteration(env, theta=0.1, discount_factor=1.0):
15     """
16     Value Iteration Algorithm. 值迭代算法
17     """
18     # 初始化状态值
19     v = np.zeros(env.nS)
20
21     # 迭代计算找到最优的状态值函数 optimal value function
22     for _ in range(50):
23         delta = 0 # 停止标志位
24
25         # 计算每个状态的状态值
26         for s in range(env.nS):
27             A = calc_action_value(s, v) # 执行一次找到当前状态的动作期望
28             best_action_value = np.max(A) # 选择最好的动作期望作为新的状态值
```

```

29
30         # 计算停止标志位
31         delta = max(delta, np.abs(best_action_value - v[s]))
32
33         # 更新状态值函数
34         v[s] = best_action_value
35
36         if delta < theta:
37             break
38
39
40     # 输出最优策略：通过最优状态值函数找到决定性策略
41     policy = np.zeros([env.ns, env.na]) # 初始化策略
42
43     for s in range(env.ns):
44         # 执行一次找到当前状态的最优状态值的动作期望 A
45         A = calc_action_value(s, v)
46
47         # 选出状态值最大的作为最优动作
48         best_action = np.argmax(A)
49         policy[s, best_action] = 1.0
50
51     return policy, v
52
53 env = gym.make("FrozenLake-v0")
54 policy, v = value_iteration(env)
55
56 print("Reshaped Policy (0=up, 1=right, 2=down, 3=left):")
57 print(np.reshape(np.argmax(policy, axis=1), [4,4]))
58 print("")
59
60 print("Reshaped Value Function:")
61 print(v.reshape([4,4]))
62 print("")

```

运行后的最终策略和值函数如下：

```

1  Reshaped Policy (0=up, 1=right, 2=down, 3=left):
2  [[0 1 2 3]
3   [0 0 0 0]
4   [1 1 0 0]
5   [0 2 1 0]]
6
7  Reshaped Value Function:
8  [[0.          0.          0.01234568 0.00411523]
9   [0.          0.          0.06995885 0.          ]
10 [0.02469136 0.14814815 0.26200274 0.          ]
11 [0.          0.3127572  0.62688615 0.          ]]

```

3.代码解析

对于给定的状态 s 计算其动作 a 的期望值

主要依据公式 $\sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$ ，通过遍历当前状态下的所有动作，获得当前状态动作的期望。

```

1  def calc_action_value(state, V, discount_factor=1.0):
2      """
3      Calculate the expected value of each action in a given state.
4      对于给定的状态 s 计算其动作 a 的期望值
5      """
6      A = np.zeros(env.nA)
7      for a in range(env.nA):
8          for prob, next_state, reward, done in env.P[state][a]:
9              A[a] += prob * (reward + discount_factor * V[next_state])
10     return A

```

在有了calc_action_value()函数后，就可以实现值迭代算法。在正式进行值迭代之前，采用np.zeros()函数将状态值向量都初始化为0。

由于折扣因子为1，可能会导致状态值的更新无法根据阈值停止的情况，因此需要使用截断的方式来控制遍历次数。在FrozenLake游戏环境中，状态的次数较少，因此迭代次数无需太多。根据经验，迭代次数设置为状态值的3-4倍（如16×3~50）就可以满足实际的迭代需求。

在算法值迭代过程中，需要将策略初始化 $[env.nS \times env.nA]$ 大小的矩阵，然后遍历每一个状态，找到使得状态值最大的动作（即最优动作），最后在策略矩阵中把该动作的位置（best action）设为1。在上述代码中，当经过50次值迭代或者满足条件 $\Delta < \theta$ 时，就认为值迭代算法已经找到最优状态 v^* ，而足有策略就是选择使得状态值最大的动作。因此，最后要做的就是根据最优状态值获得最优策略 π^* 。