

2025-2026-1 学期强化学习课程 - 第二次作业

1120231863 左逸龙

October, 29 2025

1 通勤北理工

(i) 既然我们已经知晓了最优的 Q^* 表，那么每一状态下的最优策略满足：

$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{a \in A(s)} Q^*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

因此最优策略为：

- 在状态 S_1 下，乘坐班车
- 在状态 S_2 下，乘坐班车
- 在状态 S_3 下，乘坐地铁

(ii) 此时最优策略为：

- 在状态 S_{12} 下，乘坐班车
- 在状态 S_3 下，乘坐班车

显然，此时得到的最优策略与使用真实三状态表示时得到的最优策略不同，关键区别在于原来状态 S_3 的最优策略为乘坐地铁，而现在状态 S_3 的最优策略为乘坐班车。

之所以会出现这样的变化，是因为状态聚合导致智能体 Agent 无法区分 S_1 与 S_2 的未来价值，使得决策依据退化为即时奖励。分析如下：

- Q-learning 算法的更新公式为：

$$Q(s, a) = Q(s, a) + \alpha \left[R + \gamma \max_{a' \in A(s')} Q(s', a') \right]$$

其中 α 为学习率， R 为即时奖励， γ 为折扣因子。由公式可以看出，决定 Q 值的不仅仅只有即时奖励，还有未来价值的预期。

- 在原始的三状态模型中，Agent 可以精确地知道每个动作会导向哪个具体的状态。从 S_3 出发，乘坐地铁会到达 S_2 ，而 S_2 的长期价值 ($V^*(S_2) = 1.95$) 远高于乘坐班车所到达的 S_1 的长期价值 ($V^*(S_1) = 1.65$)。尽管坐地铁的即时奖励更低，但为了追求 S_2 带来的更高未来收益，最优策略是选择乘坐地铁。

- 在聚合后的二状态模型中， S_1 和 S_2 被合并为宏状态 S_{12} 。此时，无论从 S_3 出发选择乘坐班车（到达 S_1 ）还是乘坐地铁（到达 S_2 ），在 Agent 看来，下一个状态都是同一个 S_{12} 。因此，这两个动作所带来的未来价值预期是完全相同的（都等于 $\gamma V^*(S_{12})$ ）。
- 当两个动作的未来价值预期相同时，决策的优劣就完全取决于即时奖励。根据表 1， $R(S_3, \text{班车}) = -0.5$ ，而 $R(S_3, \text{地铁}) = -0.7$ 。由于 $-0.5 > -0.7$ ，选择乘坐班车能获得更好的即时奖励。因此，在这种信息受限的情况下，最优策略从乘坐地铁转变为乘坐班车。

综上，导致这种策略上变化的原因是状态表示的粒度变粗后，Q-learning 算法泛化（或平均化）了不同状态的价值，导致决策依据从长远未来价值退化为即时奖励。

2 Frozenlake 小游戏

- (i) 以下是 `policy_evaluation`, `policy_improvement`, `policy_iteration` 的代码实现：
- `policy_evaluation`：实现了策略评估，即给定一个策略，计算其价值函数。

```

1 def policy_evaluation(
2     P: PTType,
3     nS: int,
4     nA: int,
5     policy: np.ndarray,
6     gamma: float = 0.9,
7     tol: float = 1e-3
8 ) -> np.ndarray:
9
10    value_function = np.zeros(nS)
11    while True:
12        delta = 0 # 用于记录价值函数的变化量，小于tol时认为收敛，停止迭代
13        for s in range(nS):
14            v = value_function[s]
15            a = policy[s]
16
17            new_v = 0
18
19            # 状态价值函数更新公式：
20            #  $V_{k+1}(s) = \sum_{a \in A} \pi(a|s) * (\sum_{s' \in S'} P(s'|s,a) * (R(s'|s,a) + \gamma * V_k(s')))$ 
21            # 此处由于policy已经确定，所以  $A(s) = \{a\}$ ,  $\pi(a|s) = 1$ ,
22            # 因此原式化可以简化为：
23            #  $V_{k+1}(s) = \sum_{s' \in S'} P(s'|s,a) * (R(s'|s,a) + \gamma * V_k(s'))$ 
24        for prob, next_state, reward, terminal in P[s][a]:
25            new_v += prob * (reward + gamma * value_function[next_state])
26        value_function[s] = new_v
27        delta = max(delta, abs(v - value_function[s]))
28
29    if delta < tol:
30
31

```

```
32         break
33
34     return value_function
```

- `policy_improvement` : 实现了策略提升，即给定价值函数，计算其最优策略。

```
1 def policy_improvement(
2     P: PTType,
3     nS: int,
4     nA: int,
5     value_from_policy: np.ndarray,
6     policy: np.ndarray,
7     gamma: float = 0.9
8 ) -> np.ndarray:
9
10    new_policy = np.zeros(nS, dtype="int")
11    for s in range(nS):
12        q_values = np.zeros(nA)
13        for a in range(nA):
14            # 动作价值函数更新公式:
15            #  $Q(s,a) = \sum_{s' \in S} P(s'|s,a) * (R(s'|s,a) +$ 
16            #  $\gamma V(s'))$ 
17            for prob, next_state, reward, terminal in P[s][a]:
18                q_values[a] += prob * (reward + gamma *
value_from_policy[next_state])
19
20        # 贪婪策略，选择 Q 值最大的动作
21        new_policy[s] = np.argmax(q_values)
22
23    return new_policy
```

- `policy_iteration` : 实现了策略迭代，交替进行策略评估与提升，直到收敛。

```
1 def policy_iteration(
2     P: PTType,
3     nS: int,
4     nA: int,
5     gamma: float = 0.9,
6     tol: float = 1e-3
7 ) -> Tuple[np.ndarray, np.ndarray]:
8
9     value_function = np.zeros(nS)
10    policy = np.zeros(nS, dtype=int)
11
12    iterations = 0
13    while True:
14        value_function = policy_evaluation(P, nS, nA, policy, gamma, tol)
15
16        new_policy = policy_improvement(P, nS, nA, value_function, policy,
gamma)
17
18        iterations += 1
19
20        # 如果新旧策略相同，则认为策略收敛，停止迭代
21        if np.array_equal(policy, new_policy):
22            break
23
```

```
24     policy = new_policy
25
26     print(f"Policy Iteration converged in {iterations} iterations")
27
28     return value_function, policy
```

以下是实验的关键配置，其中 `env.is_slippery = False` 表示确定性环境：

```
1 # config.yaml
2 env:
3   map_size: 4
4   frozen_prob: 0.8
5   seed: 20241022
6   is_slippery: False
7   policy_iteration:
8     gamma: 0.9
9     tol: 1e-3
10  render:
11    max_steps: 100
12  algorithm: policy_iteration
```

以下是实验结果，有修改 `run.py` 以记录更多数据：

```
1 Policy Iteration converged in 8 iterations
2 Training completed in 0.0012 seconds
3
4 Test results (100 episodes):
5 Total reward: 100.00
6 Success rate: 100.00%
7 Average steps (successful episodes): 5.00
```

可以看到，策略迭代在 8 次迭代后收敛，测试结果表明智能体在 100 次测试中成功到达目标状态 100 次，成功率为 100%，平均步数为 5 步，均为实际最优策略。可以看出策略迭代算法取得了十分理想的效果。

(ii) 以下是 Q-Learning 的代码实现：

```
1 import gymnasium
2 from datetime import datetime
3
4 def QLearning(
5     env:gymnasium.Env,
6     num_episodes=2000,
7     gamma=0.9,
8     lr=0.1,
9     epsilon=0.8,
10    epsilon_decay=0.99
11 ) -> np.ndarray:
12
13     nS:int = env.observation_space.n
14     nA:int = env.action_space.n
15     Q = np.zeros((nS, nA))
16
17     # 用于监控训练进度
```

```

18     total_rewards = []
19     success_count = []
20     start_time = datetime.now()
21
22     max_steps_per_episode = 100 # 防止无限循环
23
24     for episode in range(num_episodes):
25         state, info = env.reset()
26         done = False
27         episode_reward = 0
28         steps = 0
29
30         while not done and steps < max_steps_per_episode:
31             # 实现了 epsilon-greedy 策略:
32             # 小于 epsilon 时随机选择动作, 否则选择 Q 值最大的动作
33             if np.random.random() < epsilon:
34                 action = env.action_space.sample()
35             else:
36                 action = np.argmax(Q[state])
37
38             next_state, reward, terminated, truncated, info = env.step(action)
39             done = terminated or truncated
40             episode_reward += reward
41
42             # Q-learning 更新公式:
43             #  $Q(s,a) \leftarrow Q(s,a) + lr * [R(s'|s,a) + \gamma \max_a Q(s',a) - Q(s,a)]$ 
44             best_next_action = np.argmax(Q[next_state])
45             td_target = reward + gamma * Q[next_state, best_next_action]
46             td_error = td_target - Q[state, action]
47             Q[state, action] += lr * td_error
48
49             state = next_state
50             steps += 1
51
52
53         # 减少 epsilon, 设置下界保持一定程度的探索
54         epsilon = max(0.1, epsilon * epsilon_decay)
55
56         # 打印数据用
57         total_rewards.append(episode_reward)
58         success_count.append(1 if episode_reward > 0 else 0)
59
60         # 每 500 个 episode 打印一次进度
61         if (episode + 1) % 500 == 0:
62             end_time = datetime.now()
63             elapsed_time = (end_time - start_time).total_seconds()
64             avg_reward = np.mean(total_rewards[-500:])
65             success_rate = np.mean(success_count[-500:]) * 100
66             print(f"Episode {episode + 1}/{num_episodes}, "
67                   f"Avg Reward (last 500): {avg_reward:.3f}, "
68                   f"Success Rate (last 500): {success_rate:.3f}%, "
69                   f"Time: {elapsed_time:.2f}s, "
70                   f"Epsilon: {epsilon:.3f}")
71
72     return Q

```

以下是实验的关键配置，其中 `env.is_slippery = False` 表示确定性环境：

```
1 # config.yaml
2 env:
3   map_size: 4
4   frozen_prob: 0.8
5   seed: 20241022
6   is_slippery: False
7 qlearning:
8   num_episodes: 5000
9   gamma: 0.9
10  learning_rate: 0.1
11  epsilon: 1
12  epsilon_decay: 0.9995
13 render:
14   max_steps: 100
15 algorithm: QLearning
```

以下是实验结果，有修改 `run.py` 以记录更多数据：

```
1 Episode 500/5000, Avg Reward (last 500): 0.080, Success Rate (last 500):
2 8.000%, Time: 0.08s, Epsilon: 0.779
2 Episode 1000/5000, Avg Reward (last 500): 0.254, Success Rate (last 500):
2 25.400%, Time: 0.16s, Epsilon: 0.606
3 Episode 1500/5000, Avg Reward (last 500): 0.436, Success Rate (last 500):
3 43.600%, Time: 0.23s, Epsilon: 0.472
4 Episode 2000/5000, Avg Reward (last 500): 0.552, Success Rate (last 500):
4 55.200%, Time: 0.30s, Epsilon: 0.368
5 Episode 2500/5000, Avg Reward (last 500): 0.718, Success Rate (last 500):
5 71.800%, Time: 0.36s, Epsilon: 0.286
6 Episode 3000/5000, Avg Reward (last 500): 0.776, Success Rate (last 500):
6 77.600%, Time: 0.43s, Epsilon: 0.223
7 Episode 3500/5000, Avg Reward (last 500): 0.828, Success Rate (last 500):
7 82.800%, Time: 0.49s, Epsilon: 0.174
8 Episode 4000/5000, Avg Reward (last 500): 0.866, Success Rate (last 500):
8 86.600%, Time: 0.54s, Epsilon: 0.135
9 Episode 4500/5000, Avg Reward (last 500): 0.878, Success Rate (last 500):
9 87.800%, Time: 0.60s, Epsilon: 0.105
10 Episode 5000/5000, Avg Reward (last 500): 0.908, Success Rate (last 500):
10 90.800%, Time: 0.65s, Epsilon: 0.100
11
12 Training completed in 0.65 seconds
13
14 Test results (100 episodes):
15 Total reward: 100.00
16 Success rate: 100.00%
17 Average steps (successful episodes): 5.00
```

可以看到，Q-Learning 在训练了 0.65 秒后收敛，训练过程当中 `Avg Reward` 和 `Success Rate` 均在稳步提升，表明策略正在逐渐收敛到最优策略。测试结果表明智能体在 100 次测试中成功到达目标状态 100 次，成功率为 100%，平均步数为 5 步，均为实际最优策略。可以看出 Q-Learning 算法取得了十分理想的效果。

(iii) 保持其他配置不变，仅修改 `env.is_slippery = True`，两种算法的实验结果如下：

- Policy Iteration:

```
1 Policy Iteration converged in 5 iterations
2 Training completed in 0.0041 seconds
3
4 Test results (100 episodes):
5 Total reward: 100.00
6 Success rate: 100.00%
7 Average steps (successful episodes): 35.38
```

- Q-Learning:

```
1 Episode 500/5000, Avg Reward (last 500): 0.020, Success Rate (last 500):
2.000%, Time: 0.06s, Epsilon: 0.779
2 Episode 1000/5000, Avg Reward (last 500): 0.024, Success Rate (last 500):
2.400%, Time: 0.13s, Epsilon: 0.606
3 Episode 1500/5000, Avg Reward (last 500): 0.076, Success Rate (last 500):
7.600%, Time: 0.22s, Epsilon: 0.472
4 Episode 2000/5000, Avg Reward (last 500): 0.100, Success Rate (last 500):
10.000%, Time: 0.32s, Epsilon: 0.368
5 Episode 2500/5000, Avg Reward (last 500): 0.164, Success Rate (last 500):
16.400%, Time: 0.44s, Epsilon: 0.286
6 Episode 3000/5000, Avg Reward (last 500): 0.224, Success Rate (last 500):
22.400%, Time: 0.57s, Epsilon: 0.223
7 Episode 3500/5000, Avg Reward (last 500): 0.282, Success Rate (last 500):
28.200%, Time: 0.72s, Epsilon: 0.174
8 Episode 4000/5000, Avg Reward (last 500): 0.360, Success Rate (last 500):
36.000%, Time: 0.88s, Epsilon: 0.135
9 Episode 4500/5000, Avg Reward (last 500): 0.332, Success Rate (last 500):
33.200%, Time: 1.04s, Epsilon: 0.105
10 Episode 5000/5000, Avg Reward (last 500): 0.438, Success Rate (last 500):
43.800%, Time: 1.23s, Epsilon: 0.100
11
12 Training completed in 1.23 seconds
13
14 Test results (100 episodes):
15 The agent didn't reach a terminal state in 100 steps.
16 Total reward: 99.00
17 Success rate: 99.00%
18 Average steps (successful episodes): 38.23
```

直接观察结果，可以发现相较于确定性环境，两种算法有如下变化：

- Policy Iteration:

- 收敛速度变慢，时间从 0.0012s 增加到 0.0041s，即便迭代次数从 8 次降低到了 5 次；
- 成功率保持 100% 不变，但平均步数从 5 步增加到了 35.38 步；倘若我们提高要求，进一步限制最大步数的话，成功率可能会降低；

- Q-Learning:

- 训练时间从 0.65s 增加到 1.23s，表明随机环境需要更多采样步骤；

- 训练曲线不稳定，训练过程当中的最终成功率显著降低至 43.8%，相比确定性环境的 90.8% 下降了一半多；
- 然而测试结果的成功率依然高达 99%，这是因为测试时使用的是贪婪策略（选择 Q 值最大的动作），而训练时使用 epsilon-greedy 策略（有 10% 的探索概率），因此测试性能能够反映已学习到的最优策略，而训练过程中的成功率受探索行为影响而较低；

对比两种算法在随机性环境下的表现，可以发现：

- Policy Iteration 相较于 Q-Learning 收敛速度更快，运行时间更短，原因在于：**
 - Policy Iteration 是基于模型的（Model-based）算法，直接利用环境的转移概率矩阵 P 进行动态规划，通过迭代计算贝尔曼方程即可收敛，不需要采样；
 - Q-Learning 是无模型的（Model-free）算法，需要通过与环境大量交互来估计 Q 值，随机环境中相同的状态-动作对会产生不同的结果，需要更多样本才能准确估计；
- Policy Iteration 的鲁棒性显著优于 Q-Learning，原因在于：**
 - Policy Iteration 在策略评估阶段会计算所有可能转移的期望值，公式为 $V(s) = \sum_{s'} P(s'|s, a)[R(s, a, s') + \gamma V(s')]$ ，这种期望值计算能够很好地适应随机性环境；
 - Q-Learning 通过单次采样更新 Q 值，公式为 $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$ ，单次采样无法反映转移的真实概率分布，导致学习过程噪声大、不稳定；
 - 从实验数据看，Policy Iteration 测试成功率 100%，而 Q-Learning 测试成功率虽然达到 99%，但训练过程显示其学习曲线波动较大，最终训练成功率仅 43.8%，说明策略质量不如 Policy Iteration 稳定；

(iv) 保持其他配置不变，仅修改 `map_size = 6`，两种算法的实验结果如下：

```

1 [2025-10-28 07:30:40,403][HYDRA]      #0 : algorithm=policy_iteration
2 env.is_slippery=False env.render_mode=ansi
3 -----
4 Beginning POLICY_ITERATION
5 -----
6 Policy Iteration converged in 13 iterations
7 Training completed in 0.0056 seconds
8
9 Test results (100 episodes):
10 Total reward: 100.00
11 Success rate: 100.00%
12 Average steps (successful episodes): 9.00
13 [2025-10-28 07:30:40,520][HYDRA]      #1 : algorithm=policy_iteration
14 env.is_slippery=True env.render_mode=ansi
15 -----
16 Beginning POLICY_ITERATION
17 -----
```

```

18 Policy Iteration converged in 7 iterations
19 Training completed in 0.0081 seconds
20
21 Test results (100 episodes):
22 Total reward: 5.00
23 Success rate: 5.00%
24 Average steps (successful episodes): 415.60
25 [2025-10-28 07:30:40,649][HYDRA]      #2 : algorithm=QLearning
env.is_slippery=False env.render_mode=ansi
26
27 -----
28 Beginning QLEARNING
29 -----
30 Episode 500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.05s, Epsilon: 0.779
31 Episode 1000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.11s, Epsilon: 0.606
32 Episode 1500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.18s, Epsilon: 0.472
33 Episode 2000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.26s, Epsilon: 0.368
34 Episode 2500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.36s, Epsilon: 0.286
35 Episode 3000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.49s, Epsilon: 0.223
36 Episode 3500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.65s, Epsilon: 0.174
37 Episode 4000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.86s, Epsilon: 0.135
38 Episode 4500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 1.11s, Epsilon: 0.105
39 Episode 5000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 1.40s, Epsilon: 0.100
40
41 Training completed in 1.40 seconds
42
43 Test results (100 episodes):
44 Total reward: 0.00
45 Success rate: 0.00%
46 Average steps (successful episodes): 0.00
47 [2025-10-28 07:30:42,268][HYDRA]      #3 : algorithm=QLearning
env.is_slippery=True env.render_mode=ansi
48
49 -----
50 Beginning QLEARNING
51 -----
52 Episode 500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.06s, Epsilon: 0.779
53 Episode 1000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.10s, Epsilon: 0.606
54 Episode 1500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.14s, Epsilon: 0.472
55 Episode 2000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.18s, Epsilon: 0.368
56 Episode 2500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.21s, Epsilon: 0.286
57 Episode 3000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.24s, Epsilon: 0.223
58 Episode 3500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.27s, Epsilon: 0.174

```

```

59 Episode 4000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.30s, Epsilon: 0.135
60 Episode 4500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.33s, Epsilon: 0.105
61 Episode 5000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):
0.000%, Time: 0.36s, Epsilon: 0.100
62
63 Training completed in 0.36 seconds
64
65 Test results (100 episodes):
66 Total reward: 0.00
67 Success rate: 0.00%
68 Average steps (successful episodes): 0.00

```

保持其他配置不变，仅修改 `map_size = 8`，两种算法的实验结果如下：

```

1 [2025-10-28 07:33:13,230][HYDRA]      #0 : algorithm=policy_iteration
env.is_slippery=False env.render_mode=ansi
2
3 -----
4 Beginning POLICY_ITERATION
5 -----
6 Policy Iteration converged in 15 iterations
7 Training completed in 0.0140 seconds
8
9 Test results (100 episodes):
10 Total reward: 100.00
11 Success rate: 100.00%
12 Average steps (successful episodes): 13.00
13 [2025-10-28 07:33:13,369][HYDRA]      #1 : algorithm=policy_iteration
env.is_slippery=True env.render_mode=ansi
14
15 -----
16 Beginning POLICY_ITERATION
17 -----
18 Policy Iteration converged in 10 iterations
19 Training completed in 0.0210 seconds
20
21 Test results (100 episodes):
22 The agent didn't reach a terminal state in 100 steps.
23 The agent didn't reach a terminal state in 100 steps.
24 The agent didn't reach a terminal state in 100 steps.
25 The agent didn't reach a terminal state in 100 steps.
26 The agent didn't reach a terminal state in 100 steps.
27 The agent didn't reach a terminal state in 100 steps.
28 The agent didn't reach a terminal state in 100 steps.
29 The agent didn't reach a terminal state in 100 steps.
30 The agent didn't reach a terminal state in 100 steps.
31 The agent didn't reach a terminal state in 100 steps.
32 Total reward: 12.00
33 Success rate: 12.00%
34 Average steps (successful episodes): 303.92
35 [2025-10-28 07:33:13,549][HYDRA]      #2 : algorithm=QLearning
env.is_slippery=False env.render_mode=ansi
36
37 -----
38 Beginning QLEARNING
39 -----

```

```
40 Episode 500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.05s, Epsilon: 0.779  
41 Episode 1000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.10s, Epsilon: 0.606  
42 Episode 1500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.17s, Epsilon: 0.472  
43 Episode 2000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.26s, Epsilon: 0.368  
44 Episode 2500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.36s, Epsilon: 0.286  
45 Episode 3000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.49s, Epsilon: 0.223  
46 Episode 3500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.64s, Epsilon: 0.174  
47 Episode 4000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.84s, Epsilon: 0.135  
48 Episode 4500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 1.10s, Epsilon: 0.105  
49 Episode 5000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 1.37s, Epsilon: 0.100  
50  
51 Training completed in 1.37 seconds  
52  
53 Test results (100 episodes):  
54 Total reward: 0.00  
55 Success rate: 0.00%  
56 Average steps (successful episodes): 0.00  
57 [2025-10-28 07:33:15,156][HYDRA]      #3 : algorithm=QLearning  
env.is_slippery=True env.render_mode=ansi  
58  
59 -----  
60 Beginning QLEARNING  
61 -----  
62 Episode 500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.04s, Epsilon: 0.779  
63 Episode 1000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.08s, Epsilon: 0.606  
64 Episode 1500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.11s, Epsilon: 0.472  
65 Episode 2000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.14s, Epsilon: 0.368  
66 Episode 2500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.17s, Epsilon: 0.286  
67 Episode 3000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.21s, Epsilon: 0.223  
68 Episode 3500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.24s, Epsilon: 0.174  
69 Episode 4000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.26s, Epsilon: 0.135  
70 Episode 4500/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.29s, Epsilon: 0.105  
71 Episode 5000/5000, Avg Reward (last 500): 0.000, Success Rate (last 500):  
0.000%, Time: 0.32s, Epsilon: 0.100  
72  
73 Training completed in 0.32 seconds  
74  
75 Test results (100 episodes):  
76 Total reward: 0.00  
77 Success rate: 0.00%  
78 Average steps (successful episodes): 0.00
```

可以看到，Policy Iteration 在面对确定性环境时总是能够找到最优策略，而即便是面对随机性环境，也有机会找到策略，使智能体成功到达目标状态，但是平均步数显著增加。以上结果表明 Policy Iteration 具有较好的鲁棒性。

而 Q-Learning 在两种情况下均表现不佳，训练与测试结果均有明显异常（全 0），我推测这是因为状态数与地图边长呈平方关系，状态空间爆炸，原有的参数设置不足以使智能体充分探索状态空间。我认为可以通过调整相关的参数设置，例如增加 episodes 数量，或者调整 epsilon 等参数，来改进 Q-Learning 的训练过程，使其能够更好地适应更大的状态空间。这也反应了 Q-Learning 算法对于参数的敏感性，没有 Policy Iteration 那么鲁棒。