


# AI HW 3

Jackson Baker  
jab132@uark.edu  
011029933

October 2025

# 1 Question I

Q1

0.0	0.0	0.0	+1.00
0.0		0.0	-1.00
0.0	0.0	0.0	0.0

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} \tau(s,a,s') (R(s,a,s') + \gamma V_k(s'))$$

Actions Results  
 Up  $\rightarrow$  60% Up, 20% Left, 20% Right  
 Down  $\rightarrow$  60% Down, 20% Left, 20% Right  
 Left  $\rightarrow$  60% Left, 20% Up, 20% Down  
 Right  $\rightarrow$  60% Right, 20% Up, 20% Down

If action into wall, stay in place

$$\gamma = 0.9$$

$$V_0(s) = 0$$

$$V_1(3,3) = \left. \begin{array}{l} \text{up} : (0.6)(0 + 0.9 \times 0) + (0.2)(0 + 0.9 \times 0) + (0.2)(1 + 0.9 \times 0) = 0.2 \\ \text{down} : (0.6)(0) + (0.2)(1) + (0.2)(0) = 0.2 \\ \text{left} : (0.6)(0) + (0.2)(0) + (0.2)(0) = 0 \\ \text{right} : (0.6)(1) + (0.2)(0) + (0.2)(0) = 0.6 \end{array} \right\} V_1(3,3) = 0.6$$

First Iteration


0.0	0.0	0.6 $\rightarrow$ +1.00	
0.0		0.0	-1.00
0.0	0.0	0.0	0.0

$$V_2(3,3): \max_a \text{ is right } = (0.6)(1 + 0.9(0)) + (0.2)(0 + 0.9(0.6)) + (0.2)(0.9(0)) = (0.6)(1.0) + (0.2)(0.9(0.6)) = 0.6 + 0.2 \times 0.9 \times 0.6 = 0.708$$

$$V_2(2,3): \max_a \text{ is right } = (0.6)(0 + 0.9(0.6)) + (0.2)(0 + 0.9(0)) + (0.2)(0 + 0.9(0)) = 0.324$$

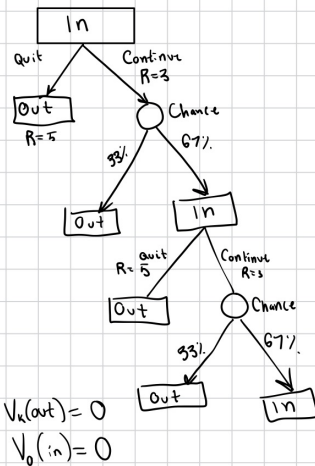
$$V_2(3,2): \max_a \text{ is up } = (0.6)(0 + 0.9(0.6)) + (0.2)(0 + 0.9(0)) + (0.2)(-1 + 0.9(0)) = 0.124$$

Second Iteration

0.0	0.324 $\rightarrow$	0.708 $\rightarrow$	+1.00
0.0		0.124 $\uparrow$	-1.00
0.0	0.0	0.0	0.0

## 2 Question II

Q2



States: In, out  
 Actions: quit, continue

Rewards:  
 $R(in, continue, out) = 3$   
 $R(in, continue, in) = 3$   
 $R(in, stop, out) = 5$

Transition Model:  
 $T(in, quit, out) = 1$   
 $T(in, cont., in) = 0.67$   
 $T(in, cont., out) = 0.33$

$$V_1(in) : \left\{ \begin{array}{l} \text{quit} : 1 \cdot (5+0) = 5 \\ \text{continue} : 0.67(3+0) + 0.33(5+0) = 3.66 \end{array} \right\} V_1(in) = 5$$


$$V_2(in) : \left\{ \begin{array}{l} \text{quit} : 5 \\ \text{continue} : 0.67(3+5) + 0.33(5+0) = 7.01 \end{array} \right\} V_2(in) = 7.01$$

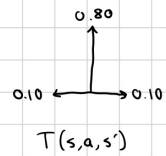
$$V_3(in) : \left\{ \begin{array}{l} \text{quit} : 5 \\ \text{continue} : 0.67(3+7.01) + 0.33(5) = 8.3567 \end{array} \right\} V_3(in) = 8.3567$$

The agent may prefer to quit because the immediate reward only given 1 iteration is larger than the reward given to continue. However, as we iterate policies, we see the agent will likely prefer to continue as soon as 2 policy iterations.

### 3 Question III

Q3

3				+1
2		Wall		-1
1				
	1	2	3	4



Q-Table

3	$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$	$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$	$\begin{matrix} 0.1 & 0.1 \\ 0 & 0 \end{matrix}$	+1
2	$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$	Wall	$\begin{matrix} -0.1 & -0.1 \\ 0 & 0.8 \end{matrix}$	-1
1	$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$	$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$	$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$	
	1	2	3	4

V-Table

3	0	0	0.8	+1
2	0	Wall	0	-1
1	0	0	0	0
	1	2	3	4

## 4 Question IV

### 4.1 Python Implementation

```
1 import numpy as np
2
3
4 class QLearningAgent:
5     def __init__(
6         self,
7         env,
8         episodes=3000,
9         alpha=0.9,
10        gamma=0.9,
11        epsilon=1.0,
12        eps_min=0.0,
13        eps_decay=0.0001,
14        max_steps=200,
15        render=False,
16        print_every=100
17    ):
18        self.env = env
19        self.episodes = episodes
20        self.alpha = alpha
21        self.gamma = gamma
22        self.epsilon = epsilon
23        self.eps_min = eps_min
24        self.eps_decay = eps_decay
25        self.max_steps = max_steps
26        self.render = render
27        self.print_every = print_every
28        self.q_table = np.zeros((env.observation_space.n,
29                                env.action_space.n))
30
31    def train(self, num_episodes=None, render_env=None):
32        self.eps_decay = 1 / (num_episodes - (0.1*num_episodes))
33        if num_episodes is None:
34            num_episodes = self.episodes
35        if render_env is None:
36            render_env = self.render
37
38        rewards = []
39        steps_per_episode = []
40
41        for ep in range(num_episodes):
42            result = self.env.reset()
43            state = result[0] if isinstance(result, tuple) else result
44
45            done, total_reward, steps = False, 0, 0
46
47            # initial Q-table and episode for rendering
48            if render_env and hasattr(self.env, 'render_mode') and
49            self.env.render_mode == "human":
50                self.env.set_episode(ep + 1)
51                self.env.set_q(self.q_table)
52                self.env.render()
```

```

51         while not done and steps < self.max_steps:
52             '''
53             ACTIONS:
54             0 = LEFT
55             1 = DOWN
56             2 = RIGHT
57             3 = UP
58             '''
59             if np.random.rand() < self.epsilon:
60                 action = self.env.action_space.sample()
61             else:
62                 action = np.argmax(self.q_table[state,:])
63
64             new_state, reward, terminated, truncated, _ =
self.env.step(action)
65             # Compute whether the episode has ended
66             done = terminated or truncated
67
68             # Calculate the TD target
69             if not terminated:
70                 best_next_action =
np.argmax(self.q_table[new_state, :])
71                 td_target = reward + self.gamma *
self.q_table[new_state, best_next_action]
72             else:
73                 td_target = reward
74
75             # Update the Q-table
76             td_error = td_target - self.q_table[state, action]
77             self.q_table[state, action] += self.alpha *
td_error
78
79             # Q-table
80             if render_env and hasattr(self.env, 'render_mode')
and self.env.render_mode == "human":
81                 self.env.set_q(self.q_table)
82
83                 self.env.render()
84
85                 state = new_state
86                 total_reward += reward
87                 steps += 1
88
89             if self.epsilon > self.eps_min:
90                 self.epsilon = max(self.epsilon - self.eps_decay,
0)
91
92             rewards.append(total_reward)
93             steps_per_episode.append(steps)
94
95             if (ep + 1) % self.print_every == 0:
96                 avg_reward = np.mean(rewards[-self.print_every:])
97                 avg_steps =
np.mean(steps_per_episode[-self.print_every:])
98                 success_rate = np.mean(
[1 if r > 0 else 0 for r in
99                 rewards[-self.print_every:]])

```

```

100         print(f"[Q] Episode {ep+1}/{num_episodes} | "
101
102         f"avg_reward({self.print_every})={avg_reward:.2f} | "
103         f"success_rate={success_rate:.2f} | "
104         f"avg_steps={avg_steps:.1f} | "
105         f"eps={self.epsilon:.3f}")
106
107     return self.q_table, rewards, steps_per_episode
108
109     def get_policy(self):
110         return np.argmax(self.q_table, axis=1)

```

## 4.2 Results for 4x4

```

[09:55:16] [~/github/ai-uark/HW3/Q4] [main *] poetry run python main.py --mode q --map 4x4 --episodes 2000
pygame 2.6.1 (SDL 2.28.4, Python 3.13.5)
Hello from the pygame community. https://www.pygame.org/contribute.html

Starting FrozenLake 4x4 - Deterministic (non-slippery) Environment

Running Q-learning...
[Q] Episode 100/2000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=7.3 | eps=0.944
[Q] Episode 200/2000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=8.2 | eps=0.889
[Q] Episode 300/2000 | avg_reward(100)=0.06 | success_rate=0.06 | avg_steps=7.6 | eps=0.833
[Q] Episode 400/2000 | avg_reward(100)=0.12 | success_rate=0.12 | avg_steps=7.9 | eps=0.778
[Q] Episode 500/2000 | avg_reward(100)=0.17 | success_rate=0.17 | avg_steps=7.7 | eps=0.722
[Q] Episode 600/2000 | avg_reward(100)=0.13 | success_rate=0.13 | avg_steps=7.4 | eps=0.667
[Q] Episode 700/2000 | avg_reward(100)=0.26 | success_rate=0.26 | avg_steps=6.9 | eps=0.611
[Q] Episode 800/2000 | avg_reward(100)=0.35 | success_rate=0.35 | avg_steps=6.9 | eps=0.556
[Q] Episode 900/2000 | avg_reward(100)=0.34 | success_rate=0.34 | avg_steps=6.4 | eps=0.500
[Q] Episode 1000/2000 | avg_reward(100)=0.42 | success_rate=0.42 | avg_steps=6.8 | eps=0.444
[Q] Episode 1100/2000 | avg_reward(100)=0.45 | success_rate=0.45 | avg_steps=6.2 | eps=0.389
[Q] Episode 1200/2000 | avg_reward(100)=0.57 | success_rate=0.57 | avg_steps=6.5 | eps=0.333
[Q] Episode 1300/2000 | avg_reward(100)=0.63 | success_rate=0.63 | avg_steps=6.8 | eps=0.278
[Q] Episode 1400/2000 | avg_reward(100)=0.80 | success_rate=0.80 | avg_steps=7.2 | eps=0.222
[Q] Episode 1500/2000 | avg_reward(100)=0.77 | success_rate=0.77 | avg_steps=6.5 | eps=0.167
[Q] Episode 1600/2000 | avg_reward(100)=0.86 | success_rate=0.86 | avg_steps=6.2 | eps=0.111
[Q] Episode 1700/2000 | avg_reward(100)=0.94 | success_rate=0.94 | avg_steps=6.2 | eps=0.056
[Q] Episode 1800/2000 | avg_reward(100)=0.99 | success_rate=0.99 | avg_steps=6.0 | eps=0.000
[Q] Episode 1900/2000 | avg_reward(100)=1.00 | success_rate=1.00 | avg_steps=6.0 | eps=0.000
[Q] Episode 2000/2000 | avg_reward(100)=1.00 | success_rate=1.00 | avg_steps=6.0 | eps=0.000
Success rate (last 100 greedy eval): 1.00
Congratulations! You've successfully solved FrozenLake!

```

### 4.3 Results for 8x8

```
[09:56:40] [~/Github/ai-uark/HW3/Q4] [main ✖] poetry run python main.py --mode q --map 8x8 --episodes 4000 --slippery
pygame 2.6.1 (SDL 2.28.4, Python 3.13.5)
Hello from the pygame community. https://www.pygame.org/contribute.html

Starting FrozenLake 8x8 - Stochastic (slippery) Environment

Running Q-learning...
[Q] Episode 100/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=32.9 | eps=0.972
[Q] Episode 200/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=33.8 | eps=0.944
[Q] Episode 300/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=36.9 | eps=0.917
[Q] Episode 400/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=37.8 | eps=0.889
[Q] Episode 500/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=34.3 | eps=0.861
[Q] Episode 600/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=37.3 | eps=0.833
[Q] Episode 700/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=36.0 | eps=0.806
[Q] Episode 800/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=30.5 | eps=0.778
[Q] Episode 900/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=32.2 | eps=0.750
[Q] Episode 1000/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=38.2 | eps=0.722
[Q] Episode 1100/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=34.2 | eps=0.694
[Q] Episode 1200/4000 | avg_reward(100)=0.02 | success_rate=0.02 | avg_steps=36.9 | eps=0.667
[Q] Episode 1300/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=32.1 | eps=0.639
[Q] Episode 1400/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=34.4 | eps=0.611
[Q] Episode 1500/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=36.6 | eps=0.583
[Q] Episode 1600/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=36.1 | eps=0.556
[Q] Episode 1700/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=35.5 | eps=0.528
[Q] Episode 1800/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=38.5 | eps=0.500
[Q] Episode 1900/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=32.8 | eps=0.472
[Q] Episode 2000/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=36.2 | eps=0.444
[Q] Episode 2100/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=33.1 | eps=0.417
[Q] Episode 2200/4000 | avg_reward(100)=0.02 | success_rate=0.02 | avg_steps=38.0 | eps=0.389
[Q] Episode 2300/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=39.9 | eps=0.361
[Q] Episode 2400/4000 | avg_reward(100)=0.00 | success_rate=0.00 | avg_steps=45.1 | eps=0.333
[Q] Episode 2500/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=39.1 | eps=0.306
[Q] Episode 2600/4000 | avg_reward(100)=0.03 | success_rate=0.03 | avg_steps=46.2 | eps=0.278
[Q] Episode 2700/4000 | avg_reward(100)=0.01 | success_rate=0.01 | avg_steps=44.8 | eps=0.250
[Q] Episode 2800/4000 | avg_reward(100)=0.05 | success_rate=0.05 | avg_steps=57.9 | eps=0.222
[Q] Episode 2900/4000 | avg_reward(100)=0.02 | success_rate=0.02 | avg_steps=51.4 | eps=0.194
[Q] Episode 3000/4000 | avg_reward(100)=0.04 | success_rate=0.04 | avg_steps=47.8 | eps=0.167
[Q] Episode 3100/4000 | avg_reward(100)=0.07 | success_rate=0.07 | avg_steps=66.0 | eps=0.139
[Q] Episode 3200/4000 | avg_reward(100)=0.07 | success_rate=0.07 | avg_steps=54.3 | eps=0.111
[Q] Episode 3300/4000 | avg_reward(100)=0.11 | success_rate=0.11 | avg_steps=64.4 | eps=0.083
[Q] Episode 3400/4000 | avg_reward(100)=0.10 | success_rate=0.10 | avg_steps=68.1 | eps=0.056
[Q] Episode 3500/4000 | avg_reward(100)=0.23 | success_rate=0.23 | avg_steps=83.7 | eps=0.028
[Q] Episode 3600/4000 | avg_reward(100)=0.30 | success_rate=0.30 | avg_steps=94.0 | eps=0.000
[Q] Episode 3700/4000 | avg_reward(100)=0.55 | success_rate=0.55 | avg_steps=105.3 | eps=0.000
[Q] Episode 3800/4000 | avg_reward(100)=0.37 | success_rate=0.37 | avg_steps=92.5 | eps=0.000
[Q] Episode 3900/4000 | avg_reward(100)=0.56 | success_rate=0.56 | avg_steps=107.7 | eps=0.000
[Q] Episode 4000/4000 | avg_reward(100)=0.51 | success_rate=0.51 | avg_steps=95.6 | eps=0.000
Success rate (last 100 greedy eval): 0.74
Congratulations! You've successfully solved FrozenLake!
```



## 4.4 Commentary on results

In the 4x4 deterministic configuration, learning is rapid and achieves a 100 percent success rate. Starting at 0 percent success at episode 100, the agent initially takes random exploratory moves until it discovers the reward. As epsilon decays, random moves decrease while exploitation of learned knowledge increases. By episode 1800, epsilon reaches 0, meaning no further exploration occurs—only optimal moves are executed. This results in a 100 percent success rate because the deterministic environment allows the agent to reliably repeat a working solution every time. The 8x8 configuration with slippery mode presents a different challenge. The possibility of slipping necessitated modifications to the default hyperparameters: alpha (learning rate), gamma (discount factor), and epsilon (exploration rate). The most impactful change was modifying epsilon decay to a linear function with an x-intercept at 10 percent of the total episode count. This approach relies on achieving early successes; as shown in the results, when the agent finds a reward within the first 1000 episodes, it establishes a foundation of successful moves to build upon. With initially high epsilon enabling sufficient exploration, the agent reliably achieves over 50 percent success rate by the final episodes.

Regarding optimality, the 4x4 deterministic policy is excellent, achieving 100 percent success in exactly 6 steps—the shortest possible path. However, the 8x8 stochastic configuration presents greater challenges. Achieving any success above 0 percent heavily depends on discovering rewards early during high-exploration phases. With linear epsilon decay, success rates spike dramatically near the end as the agent increasingly exploits the best moves learned during early random exploration. Despite this improvement, I do not believe the policy is optimal for the 8x8 stochastic configuration, likely due to insufficient training episodes. More episodes could allow the agent to discover safer paths where probability favors success and hole-avoidance is more likely. While no path may guarantee 100 percent success in a stochastic environment, I believe better solutions exist beyond the observed 51-74 percent success rate.

The fundamental difference between deterministic and stochastic learning is convergence reliability. In deterministic problems, convergence is guaranteed. Once the agent learns an optimal solution, executing the same action sequence produces identical outcomes and rewards every time. In stochastic problems, convergence is not guaranteed. A path that succeeds in one episode may fail in another due to random state transitions, making it inherently impossible to achieve perfect reliability regardless of policy quality. This uncertainty requires different learning strategies, including potentially extended exploration and acceptance of suboptimal success rates.