

```
[1]: import sys
import time
from IPython.display import Image
import gym
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.animation as animation
from frozen_lake import MDP

In [2]: def simple_plot(
    arr,
    xlabel: str,
    ylabel: str,
    title: str,
    save_path: str = "",
    show: bool = True,
    timeout: int = None,
) -> None:
    """
    if len(arr) == 2:
        plt.plot(arr[0], arr[1])
    else:
        plt.plot(arr)
    plt.xlabel(xlabel)
    plt.ylabel(ylabel)
    plt.title(title)
    if save_path:
        plt.savefig(save_path)
    if timeout is not None:
        plt.pause(timeout)
        plt.draw()
        plt.waitforbuttonpress(timeout=5)
        plt.close()
    elif show:
        plt.show()
    plt.close()

In [3]: # Create the environment
env = gym.make("FrozenLake-v1")
env.reset()

The surface is described using a grid like the following:
SFFF      (S: starting point, safe)
FFFH      (F: frozen surface, safe)
HFFF      (H: hole, fall to your doom)
GFFF      (G: goal, where the frisbee is located)
print("Observation space: ", env.observation_space)
print("Action space: ", env.action_space)

Observation space: Discrete(16)
Action space: Discrete(4)

In [4]: # Create MDP from the env as a reference
mdp = MDP(env)

In [5]: actions = [
    "left", 0,
    "down", 1,
    "right", 2,
    "up", 3,
]

act_seq = (2 * ("Right")) + (3 * ("Down")) + ("Right")
print(f"Action sequence: {act_seq}")

env.render()

for a in act_seq:
    obs, rew, done, info = env.step(actions[a])
    env.render()
    print(f"Reward: {rew}, {a}")
    print(info, "\n")
    if done:
        break

Action sequence: ("Right", "Right", "Down", "Down", "Down", "Right")

SFFF
FFFH
FFFH
FFFG
Reward: 0.00
('prob': 0.3333333333333333)

('Right')
SFFF
FFFH
FFFH
FFFG
Reward: 0.00
('prob': 0.3333333333333333)

('Down')
SFFF
FFFH
FFFH
FFFG
Reward: 0.00
('prob': 0.3333333333333333)

('Down')
SFFF
FFFH
FFFH
FFFG
Reward: 0.00
('prob': 0.3333333333333333)

('Down')
SFFF
FFFH
FFFH
FFFG
Reward: 0.00
('prob': 0.3333333333333333)

('Right')
SFFF
FFFH
FFFH
FFFG
Reward: 0.00
('prob': 0.3333333333333333)
```

Describe the environment state and action spaces, and reward function. Given a state and an action, is the state transition deterministic?

State space: $S \in \{0, 1, 2, \dots, 15\}$ represents the index from top-left to bottom-right of 4x4 grid. Each state s can be $\{S, F, H, G\}$ where

- S: starting point, safe
- F: frozen surface, safe
- H: hole, fall to your doom
- G: goal, where the frisbee is located

The terminal state is the goal state G and hole state H .

Reward is 0 for every step taken in $\{S, F, H\}$, 1 for reaching the final goal state G .

$$r = 1 \text{ if } s = G$$
$$r = 0 \text{ Otherwise}$$

The state transition is not deterministic, because the transition probability of given state and action is not 1.

Starting with the definition of a value function, show that for a deterministic policy $\pi(s)$, the value function $v(s)$ can be expressed as:

$$v(s) = \sum_{a \in S} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$$

Return $G_t = \sum_{t'=t}^T \gamma^{t'-t} R_{t'}$

Assume probabilistic transitions $T(s, a, s') = R(s|a, s)$

Deterministic policy $\pi(s) = a \forall a \in A$

$$v(s) = E_{\pi} [G_t | S_t = s]$$
$$v(s) = E_{\pi} \left[\sum_{k=t}^{\infty} \gamma^{k-t} R_k | s_t = s \right]$$
$$= E_{\pi} [R_t + \gamma \sum_{k=t+1}^{\infty} \gamma^{k-t-1} R_k | s_{t+1} = s']$$
$$= \sum_a \pi(a|s) \sum_{s'} p(s'|a, s) [r(s, a, s') + \gamma E_{\pi} \left[\sum_{k=t+1}^{\infty} \gamma^{k-t-1} R_k | s_{t+1} = s' \right]]$$
$$= \sum_a \pi(a|s) \sum_{s'} p(s'|a, s) [r(s, a, s') + \gamma v(s')]$$

In our case, $a = \pi(s)$ is a deterministic policy, we can omit the probability of policy in the above equation.

$$v(s) = \sum_{a \in S} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$$

Write a function `TestPolicy(policy)`, that returns the average rate of successful episodes over 100 trials for a deterministic policy. What is the success rate of a policy (number of times completed / total number of trials) given by $\pi(s) = (s + 1)$.

```
[6]: def TestPolicy(
    self,
    policy: Callable,
    trials: int = 100,
    render: bool = False,
    verbose: bool = False,
) -> float:
    """
    Test a policy by running it in the given environment.

    :param policy: A policy to run.
    :param env: The environment to run the policy in.
    :param render: Whether to render the environment.
    :returns: success rate over # of trials.
    """
    assert trials > 0 and isinstance(trials, int)

    success = 0
    reward = 0
    for _ in range(trials):
        obs = self.env.reset()
        done = False
        while not done:
            act = policy(obs)
            obs, rew, done, info = self.env.step(act)
            reward += rew
            if render:
                self.env.render()
            time.sleep(0.1)
            if done and obs == 15:
                success += 1
        success_rate = success / trials
        mean_reward = reward / trials
        if verbose:
            print(f"Success rates: {success_rate}")
    return success_rate, mean_reward

In [7]: # 3. Naive policy
policy = lambda s: (s + 1) % 4
naive_success_rates = []
for _ in range(10):
    naive_success_rate, _ = mdp.TestPolicy(policy, render=False)
    naive_success_rates.append(naive_success_rate)
print(f"Average naive_success_rates: {np.mean(naive_success_rates)}")

Average naive_success_rates: 0.017
```

Write a function `LearnModel`, that returns the transition probabilities $p(s'|a, s)$ and reward function $r(s, a, s')$. Estimate these values over 10^6 random samples.

```
[8]: def learnModel(self, n_samples: int = 10 ** 5) -> Tuple[np.ndarray, np.ndarray]:
    """
    Estimate transition probabilities  $p(s'|a, s)$  and
    reward function  $r(s, a, s')$  over  $n\_samples$  random samples.

    :returns: transition probabilities and reward function.
    """
    assert n_samples > 0 and isinstance(n_samples, int)

    # Dimension of observation space and action space (both discrete)
    P = np.zeros((self.n_s, self.n_s)) # transition probability:  $S \times A \times S' \rightarrow [0, 1]$ 
    R = np.zeros_like(P) # reward  $r(s, a, s')$ 

    obs = self.env.reset()
    done = False
    for _ in range(n_samples):
        # Random action:  $a \sim \pi(a|s)$ 
        act = env.action_space.sample()
        nxt_obs, rew, done, _ = self.env.step(act)

        P[obs, act, nxt_obs] += 1
        R[obs, act, nxt_obs] += rew
        obs = nxt_obs

    if done:
        # Normalize transition probabilities -> [0,1]
        pi = P.copy() # Don't modify P directly
        for a in range(self.n_a):
            total_counts = np.sum(P[:, a, :])
            pi[:, a, :] /= total_counts
        # Avoid division by zero error
        R[np.where(P == 0, R, 1)]

    # Store the estimated transition probabilities and reward function
    self.P_hat = P
    self.R_hat = R
    return pi, R
```

```
[9]: P, R_hat = mdp.learnModel(n_samples=10 ** 5)
mse = lambda x, y: np.mean((x - y) ** 2)

MSE_P, MSE_R = mse(mdp.P, P) # mse(mdp.R, R_hat)
print(f"Mean square error of P: {MSE_P}, {np.mean(MSE_R)}")

Mean square error of P: 0.019583992668908274,
Mean square error of R: 0.0013020833333333333
```

Write a function `PolicyEval()` for evaluating a given deterministic policy and with the help of this function implement a policy iteration method to solve this environment over 50 iterations. Plot the average rate of success of the learned policy at every iteration.

```
[10]: def PolicyEval(
    self,
    V: np.ndarray,
    policy: np.ndarray,
    gamma: float,
    theta: float,
) -> float:
    """
    Policy evaluation

    :param V: value function
    :param policy: a policy
    :param gamma: discount factor
    :param theta: tolerance or termination threshold
    """
    assert 0 < gamma < 1
    assert 0 < theta <= 1e-2, "Threshold should be a small positive number"

    # Using the estimations of the transition probabilities and reward function
    if self.P_hat is None or self.R_hat is None:
        self.learnModel()

    while True:
        delta = 0.0
        for s in range(self.n_s):
            act = policy[s]
            V_old = V[s]
            Vs = self.P_hat[s, act, nxt_s] * (self.R_hat[s, act, nxt_s] + gamma * V[nxt_s])
            # calculate delta
            delta = max(delta, abs(Vs - V[s]))
            # Update V
            V[s] = Vs
        if delta < theta:
            break
    return V
```

```
[11]: def PolicyIteration(
    self,
    max_iter: int = 50,
    gamma: float = 0.99,
    theta: float = 1e-8,
) -> float:
    """
    Policy iteration

    :param policy: a policy
    :param max_iter: maximum number of iterations
    :param gamma: discount factor
    :param theta: tolerance or termination threshold
    """
    assert max_iter > 0 and isinstance(max_iter, int)
    assert 0 < gamma < 1
    assert 0 < theta <= 1e-2, "Theta should be a small positive number"

    # Initialize V(s), V(pi(s))
    V = np.zeros(self.n_s)
    P_hat = np.zeros(self.n_s, dtype=int) # since actions are integers
    success_rates = []
    mean_rewards = []

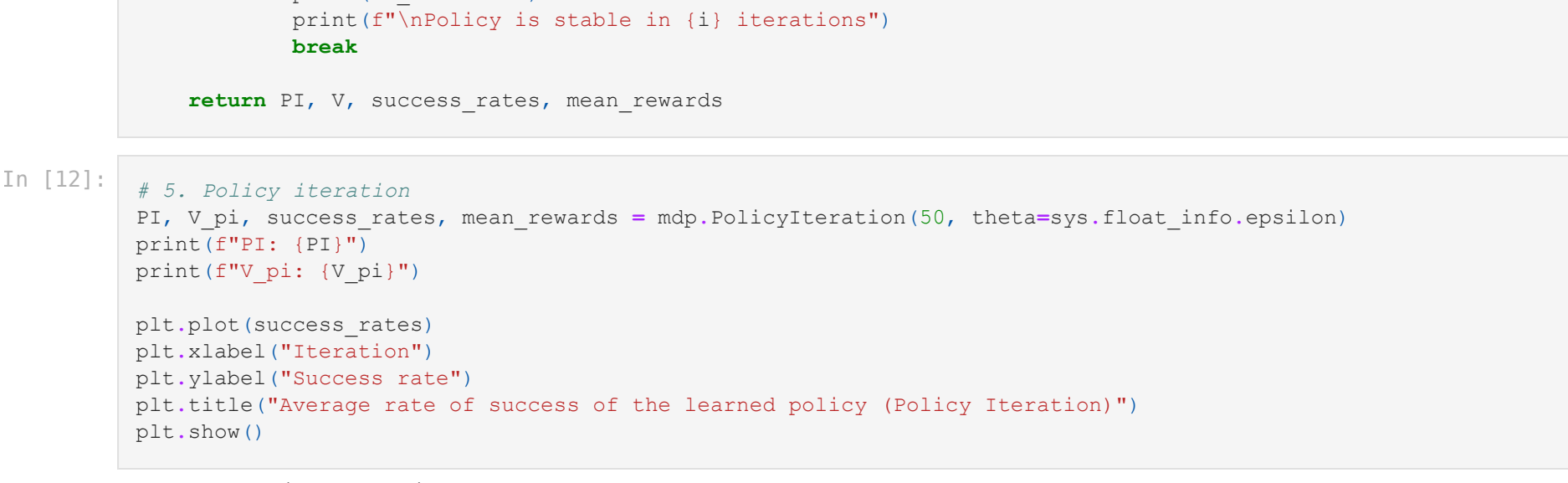
    print(f"\n----- Policy Iteration -----")
    for i in range(max_iter):
        P_hat = P.copy()
        print(f"Iteration {i+1}: ", end='')
        # Policy Evaluation
        V = self.PolicyEval(V, P_hat, gamma, theta)
        # Policy Improvement
        P_hat = self.PolicyImprovement(V, gamma)

        P_hat_fn = lambda s: P_hat[s]
        success_rate, mean_reward = self.TestPolicy(P_hat_fn, trials=100, render=False, verbose=True)
        success_rates.append(success_rate)
        mean_rewards.append(mean_reward)

        if np.all(P_hat == P_hat_fn):
            print(f"\nPolicy is stable in {i} iterations")
            break
    return P_hat, V, success_rates, mean_rewards
```

```
[12]: # 5. Policy iteration
P_hat, V_hat, success_rates, mean_rewards = mdp.PolicyIteration(50, theta=sys.float_info.epsilon)
print(f"P_hat: {P_hat}")
print(f"V_hat: {V_hat}")

plt.plot(success_rates)
plt.xlabel("Iteration")
plt.ylabel("Success rate")
plt.title("Average rate of success of the learned policy (Policy Iteration)")
plt.show()
```



Write a function `ValueIter()` that returns a deterministic policy learned through value-iteration over 50 iterations. Plot the average rate of success of the learned policy at every iteration.

```
[13]: def ValueIter(
    self,
    max_iter: int = 50,
    gamma: float = 0.99,
    theta: float = 1e-8,
) -> float:
    """
    Value iteration

    :param max_iter: maximum number of iterations
    :param gamma: discount factor
    :param theta: tolerance or termination threshold
    """
    assert max_iter > 0 and isinstance(max_iter, int)
    assert 0 < gamma < 1

    # Initialize V(s), V(pi(s))
    V = np.zeros(self.n_s)
    P_hat = np.zeros(self.n_s, dtype=int) # since actions are integers
    success_rates = []
    mean_rewards = []

    print(f"\n----- Value Iteration -----")
    for i in range(max_iter):
        print(f"Iteration {i+1}: ", end='')
        # Bellman optimality
        V_old = V[s]
        for s in range(self.n_s):
            # V(s) = max_a Q(s, a)
            Q = self.QValue(V, s, gamma) # Q(s, a, s') -> Vector of Q-values
            V[s] = max(Q)

            delta = max(delta, abs(V[s] - V_old))
        if delta < theta:
            break

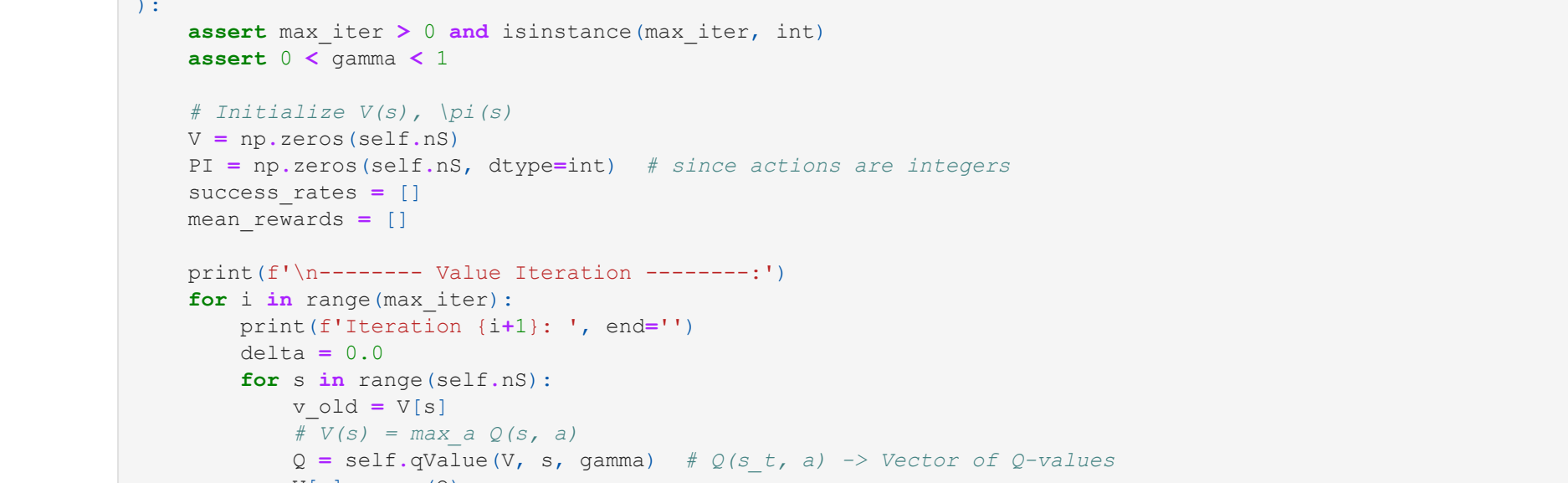
    P_hat = self.PolicyImprovement(V, gamma)

    P_hat_fn = lambda s: P_hat[s]
    success_rate, mean_reward = self.TestPolicy(P_hat_fn, trials=100, render=False, verbose=True)
    success_rates.append(success_rate)
    mean_rewards.append(mean_reward)

    return P_hat, V, success_rates, mean_rewards
```

```
[14]: # 6. Value iteration
P_hat, V_hat, success_rates, mean_rewards = mdp.ValueIter(50, theta=sys.float_info.epsilon)
print(f"P_hat: {P_hat}")
print(f"V_hat: {V_hat}")

plt.plot(success_rates)
plt.xlabel("Iteration")
plt.ylabel("Success rate")
plt.title("Average rate of success of the learned policy (Value Iteration)")
plt.show()
```



Solve the environment using Q-learning over 5000 episodes. For exploration during training, take random actions with probability $1/\epsilon$ 5000 where ϵ is the number of current episode. Plot the success rate of the learned policy at an interval of 100 episodes.

(a) Train the policy using the following learning rates with $\gamma = 0.99$. Report what you observe.

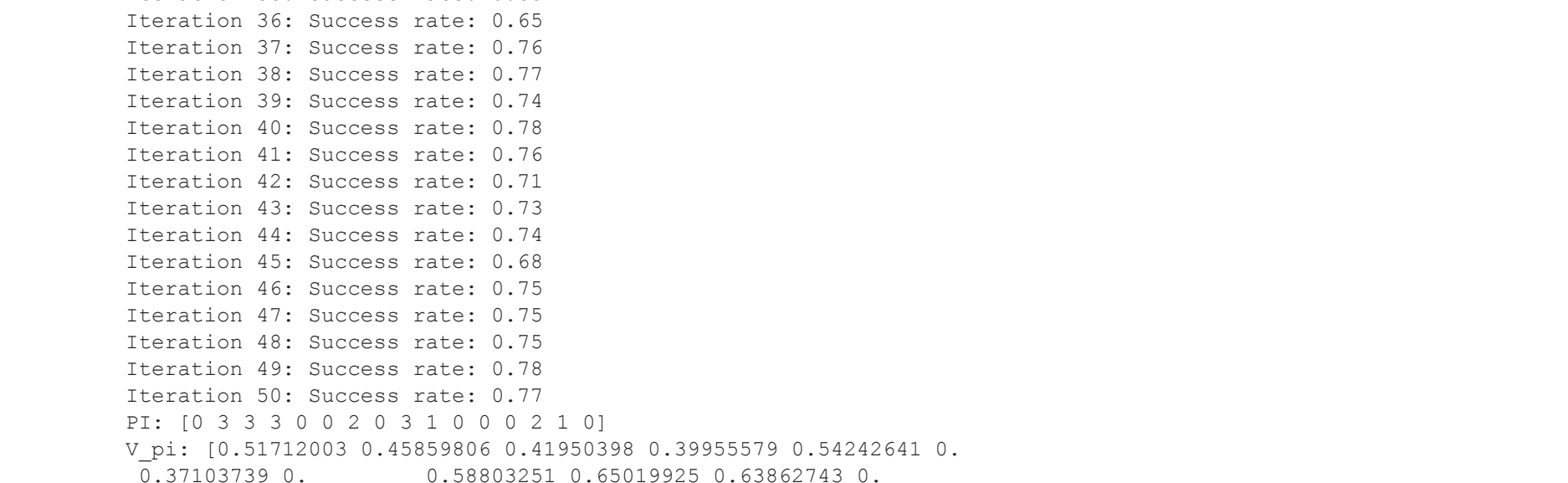
A small α tends to work better than

With a larger alpha, the Q Learning will add more weight to the reward to go. However, Q-Learning learn slowly. The Q value is not very good during the early stage of learning, (agent likely visited every state in the environment)

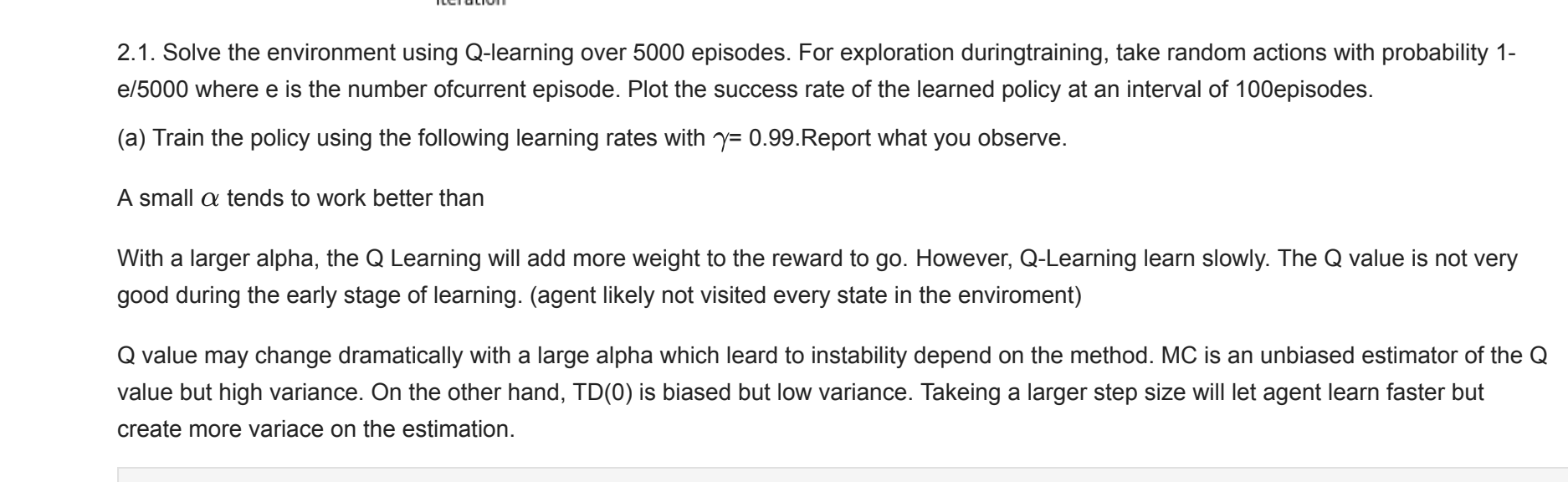
Q value may change dramatically with a large alpha which lead to instability depend on the method. MC is an unbiased estimator of the Q value but high variance. On the other hand, TD(0) is biased but low variance. Taking a larger step size will let agent learn faster but create more variance on the estimation.

```
[15]: # 7
learning_rate = [0.05, 0.1, 0.25, 0.5]
discount_factor = [0.5, 0.95, 0.99]
for lr in learning_rate:
    for QL, Q_gamma, success_rates_QL = mdp.QLearning(
        n_episodes=5_000,
        gamma=Q_gamma,
        alpha=lr,
        strategy="epsilon",
        verbose=False,
    ):
        simple_plot(
            success_rates_QL,
            xlabel="Iteration",
            ylabel="Success rate",
            title=f"Average rate of success of the learned policy (Q Learning), "
                + f"alpha={lr}"
                + f"gamma={Q_gamma}"
                + f"discount={discount_factor}"
                + f"epsilon={lr}"
        )
        # save path=plot_dir / f'QL_{lr}_{Q_gamma}.png'
        show=True,
    )
    print("\n")

----- Q Learning (alpha=0.05, gamma=0.99, strategy: epsilon) -----
Average rate of success of the learned policy (Q Learning), alpha=0.05, gamma=0.99
```



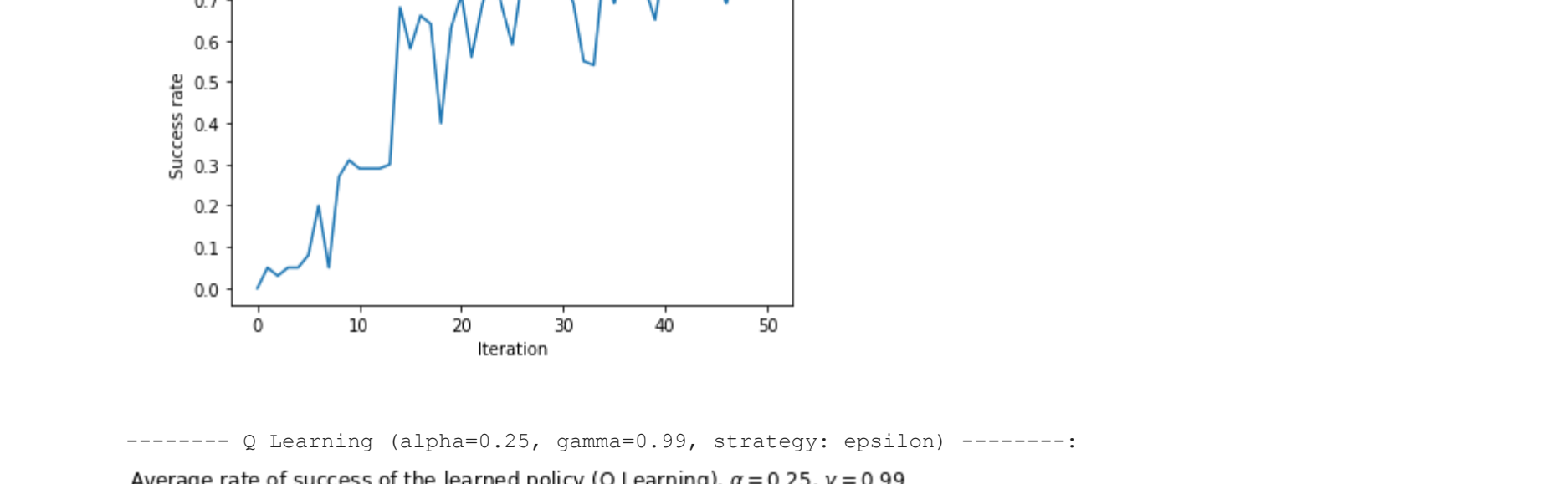
```
----- Q Learning (alpha=0.1, gamma=0.99, strategy: epsilon) -----
Average rate of success of the learned policy (Q Learning), alpha=0.1, gamma=0.99
```



```
----- Q Learning (alpha=0.25, gamma=0.99, strategy: epsilon) -----
Average rate of success of the learned policy (Q Learning), alpha=0.25, gamma=0.99
```



```
----- Q Learning (alpha=0.5, gamma=0.99, strategy: epsilon) -----
Average rate of success of the learned policy (Q Learning), alpha=0.5, gamma=0.99
```

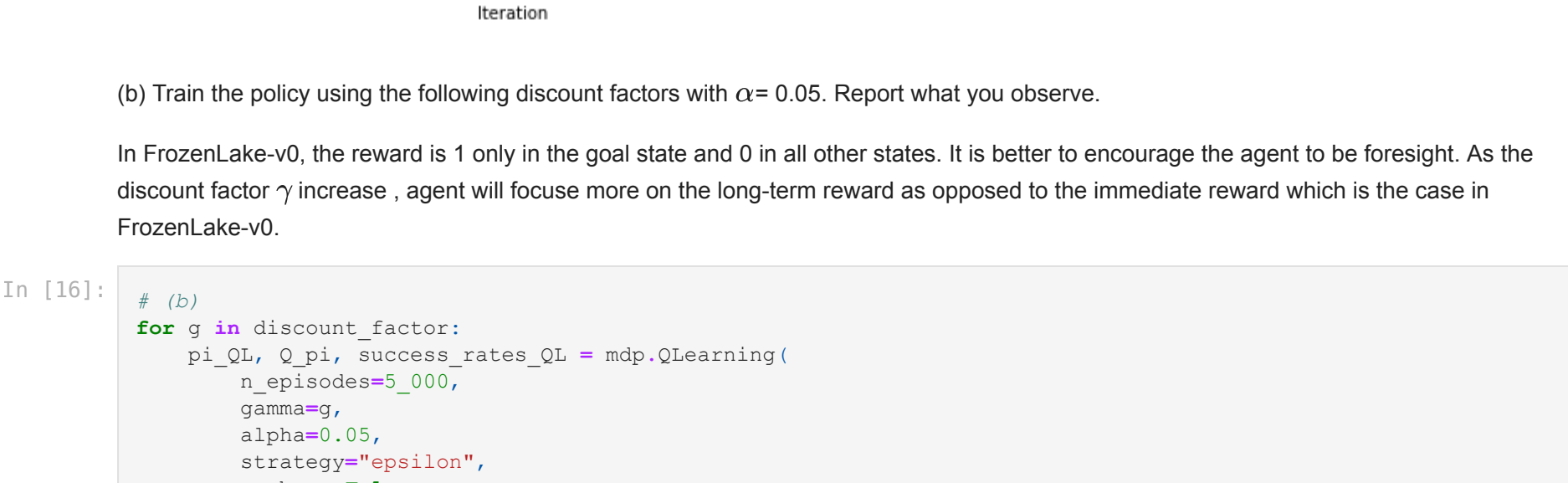


(b) Train the policy using the following discount factors with $\alpha = 0.05$. Report what you observe.

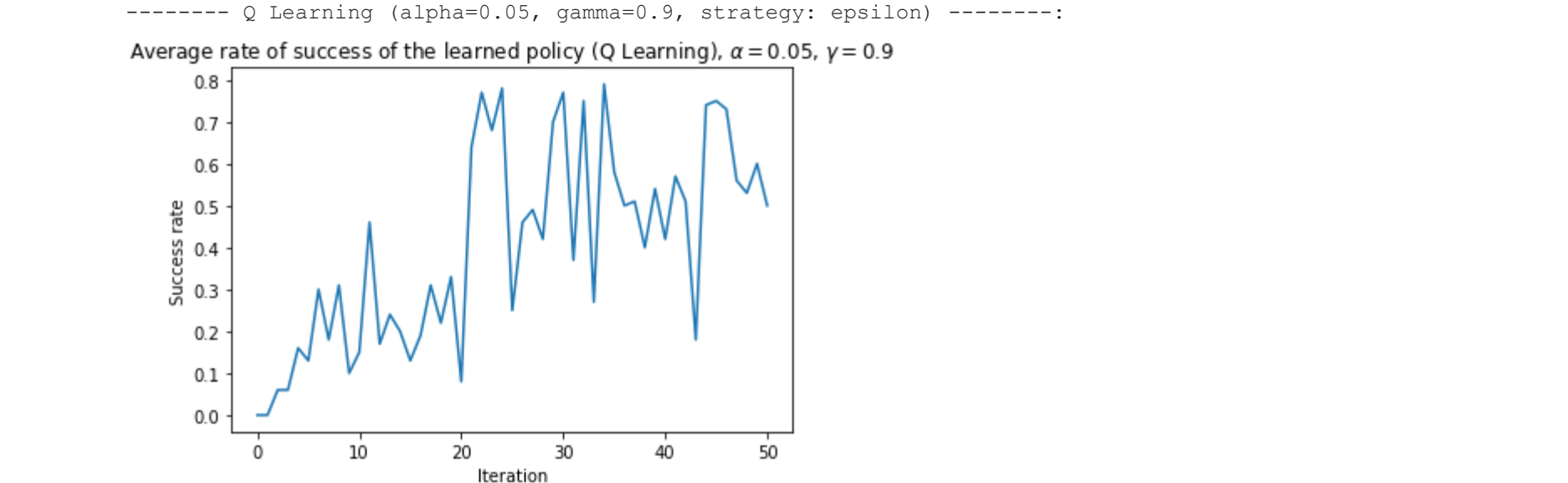
In FrozenLake-v0, the reward is 1 only in the goal state and 0 in all other states. It is better to encourage the agent to be foresight. As the discount factor γ increase, agent will focus more on the long-term reward as opposed to the immediate reward which is the case in FrozenLake-v0.

```
[16]: # (b)
for QL, Q_gamma, success_rates_QL = mdp.QLearning(
    n_episodes=5_000,
    gamma=Q_gamma,
    alpha=0.05,
    strategy="epsilon",
    verbose=False,
):
    simple_plot(
        success_rates_QL,
        xlabel="Iteration",
        ylabel="Success rate",
        title=f"Average rate of success of the learned policy (Q Learning), "
            + f"alpha={alpha}"
            + f"gamma={gamma}"
            + f"discount={discount_factor}"
            + f"epsilon={epsilon}"
        )
    # save path=plot_dir / f'QL_{alpha}_{Q_gamma}_{epsilon}.png'
    show=True,
    )
    print("\n")

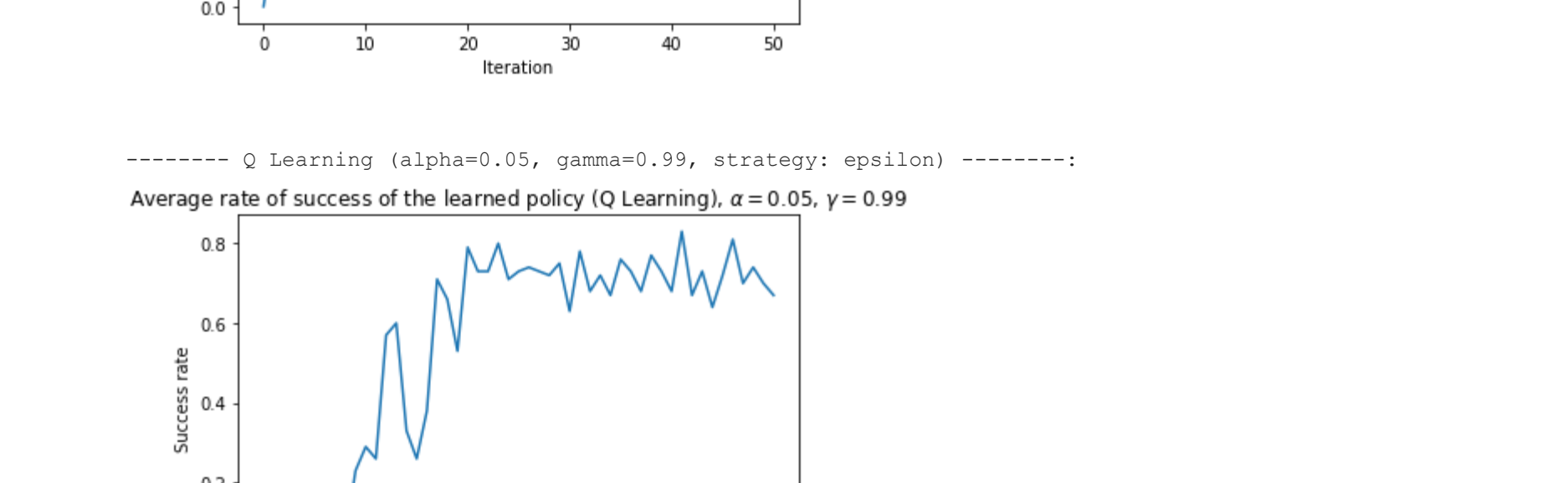
----- Q Learning (alpha=0.05, gamma=0.9, strategy: epsilon) -----
Average rate of success of the learned policy (Q Learning), alpha=0.05, gamma=0.9
```



```
----- Q Learning (alpha=0.05, gamma=0.95, strategy: epsilon) -----
Average rate of success of the learned policy (Q Learning), alpha=0.05, gamma=0.95
```



```
----- Q Learning (alpha=0.05, gamma=0.99, strategy: epsilon) -----
Average rate of success of the learned policy (Q Learning), alpha=0.05, gamma=0.99
```

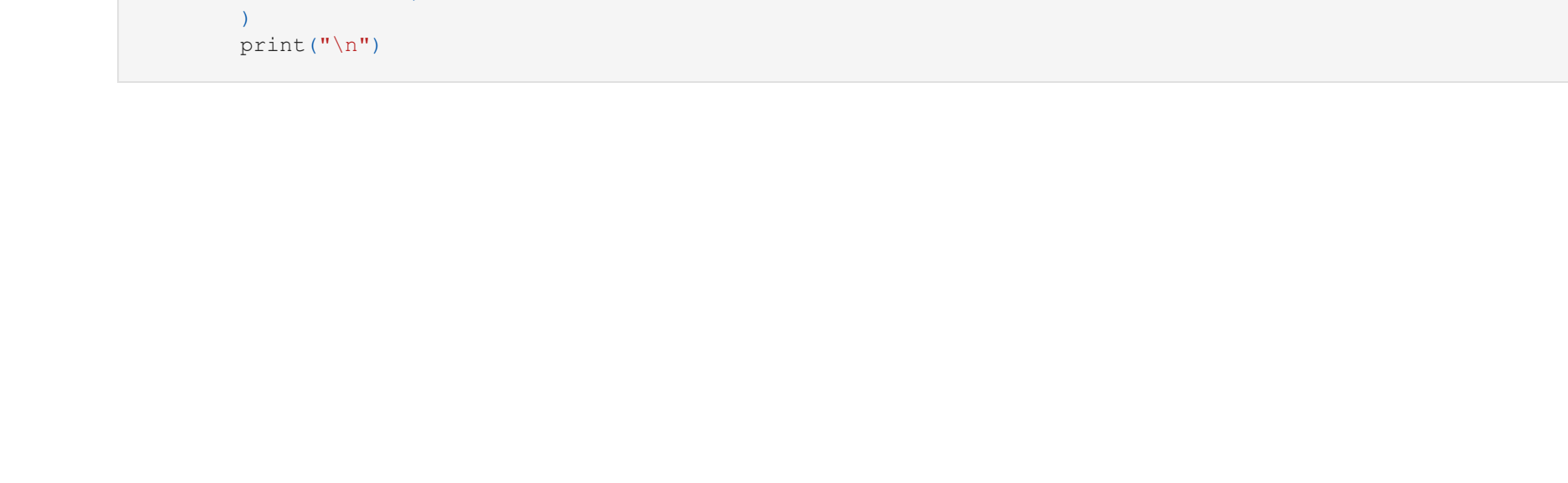


In the previous question, the exploration was linearly annealed. Solve the environment using Q-learning by proposing a different strategy to explore. Find a suitable α and γ for your method. Report your strategy and training results.

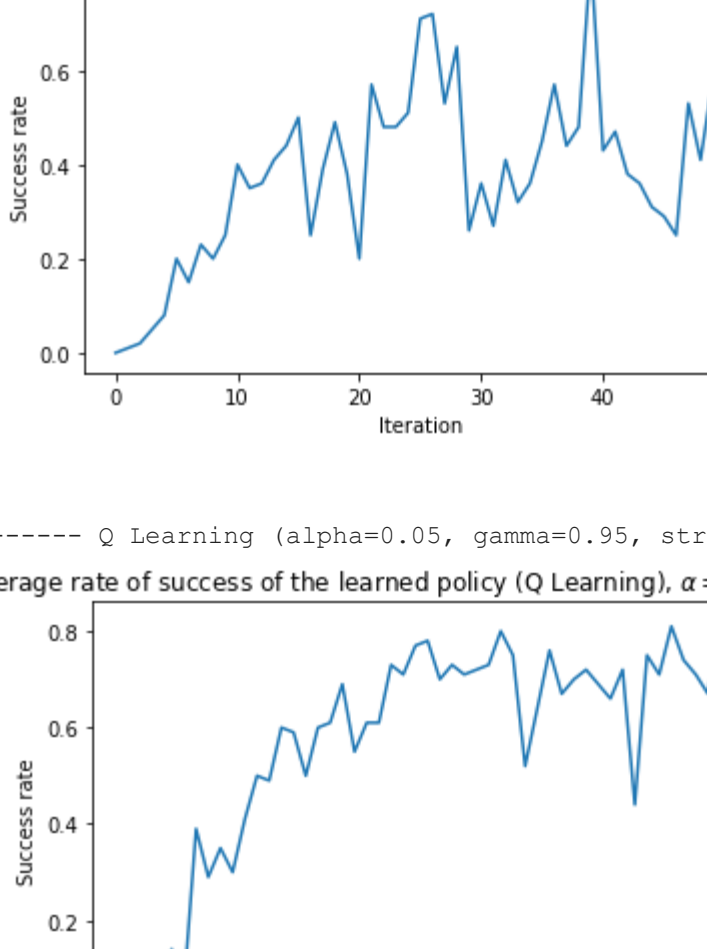
Best $\alpha = 0.1, \gamma = 0.99$

```
[17]: for lr in learning_rate:
    for QL, Q_gamma, success_rates_QL = mdp.QLearning(
        n_episodes=5_000,
        gamma=Q_gamma,
        alpha=Q_gamma,
        strategy="exponential", # "epsilon"
        verbose=False,
    ):
        simple_plot(
            success_rates_QL,
            xlabel="Iteration",
            ylabel="Success rate",
            title=f"Average rate of success of the learned policy (Q Learning), "
                + f"alpha={alpha}"
                + f"gamma={gamma}"
                + f"discount={discount_factor}"
                + f"epsilon={epsilon}"
        )
        # save path=plot_dir / f'QL_exp_{alpha}_{Q_gamma}_{epsilon}.png'
        show=True,
    )
    print("\n")

----- Q Learning (alpha=0.05, gamma=0.99, strategy: epsilon) -----
Average rate of success of the learned policy (Q Learning), alpha=0.05, gamma=0.99
```

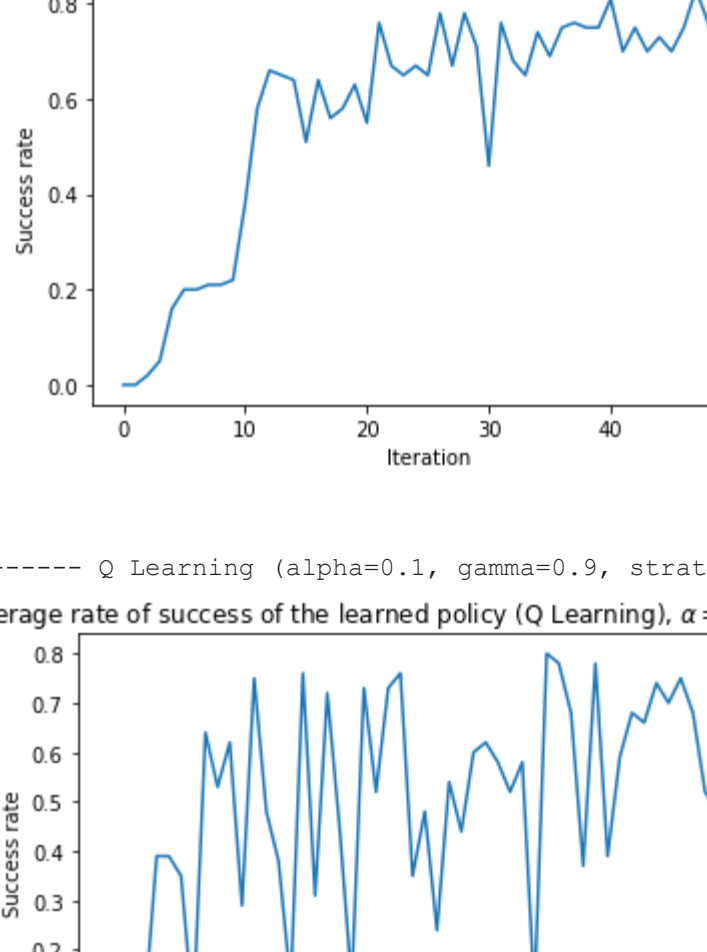


Average rate of success of the learned policy (Q Learning). $\alpha=0.05, \gamma=0.9$



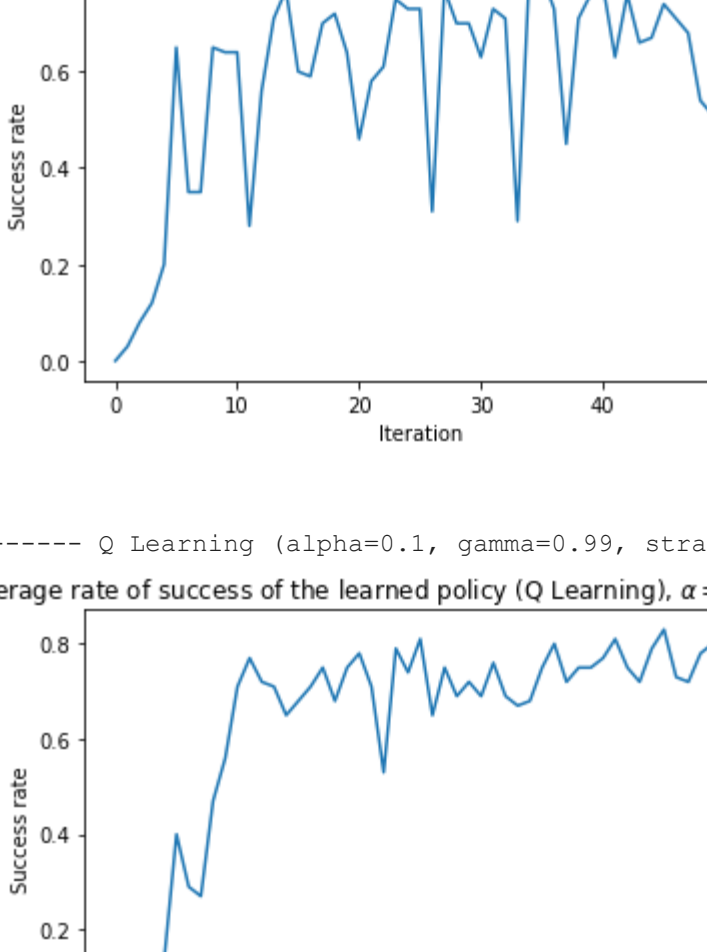
----- Q Learning (alpha=0.05, gamma=0.95, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.05, \gamma=0.95$



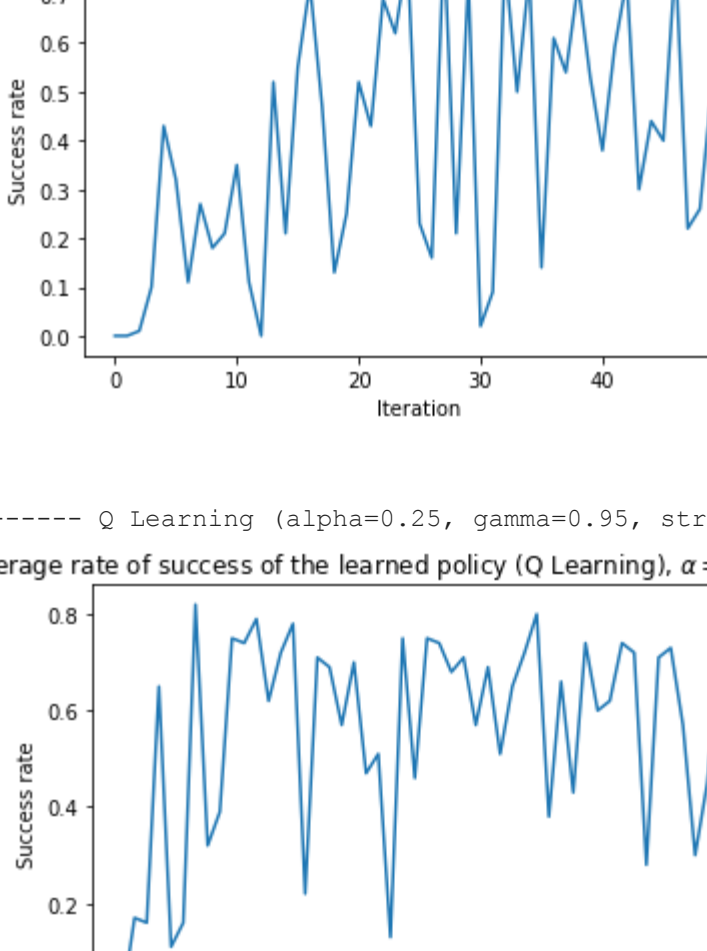
----- Q Learning (alpha=0.05, gamma=0.99, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.05, \gamma=0.99$



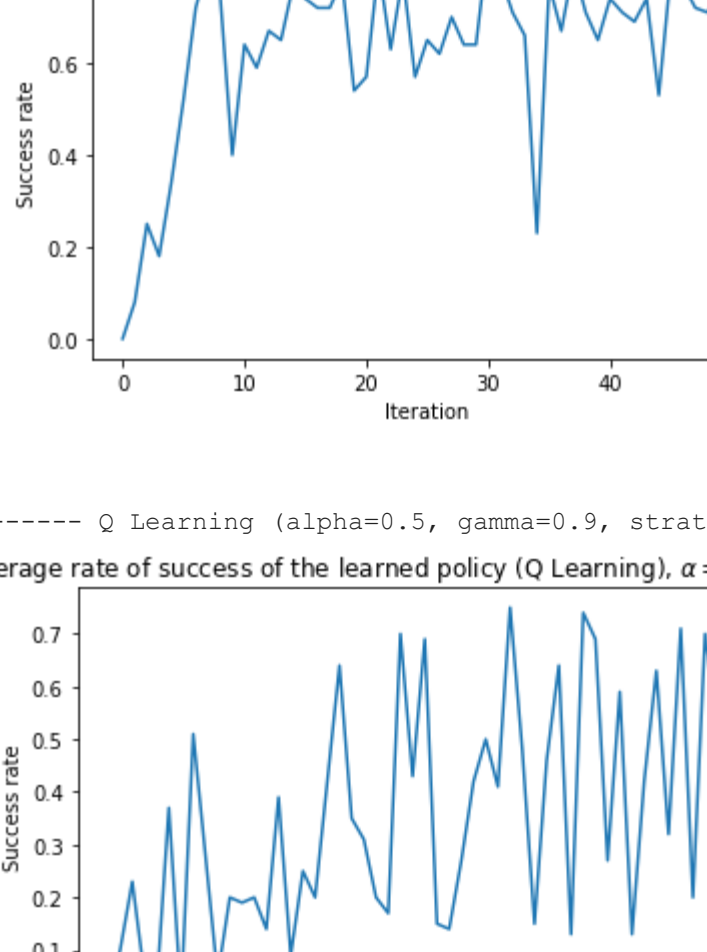
----- Q Learning (alpha=0.1, gamma=0.9, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.1, \gamma=0.9$



----- Q Learning (alpha=0.1, gamma=0.95, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.1, \gamma=0.95$



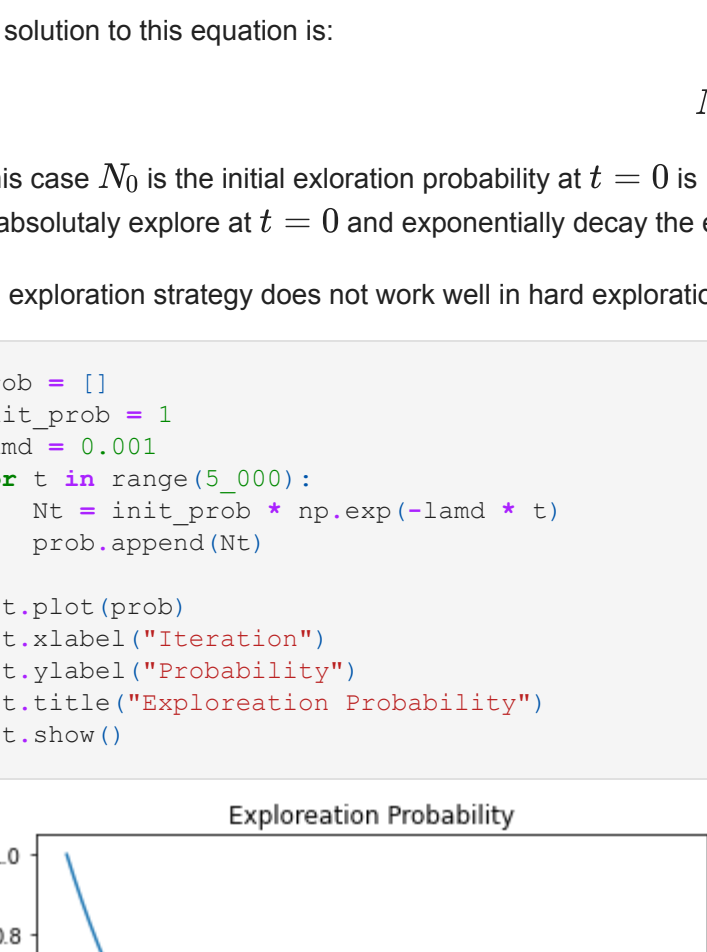
----- Q Learning (alpha=0.1, gamma=0.99, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.1, \gamma=0.99$



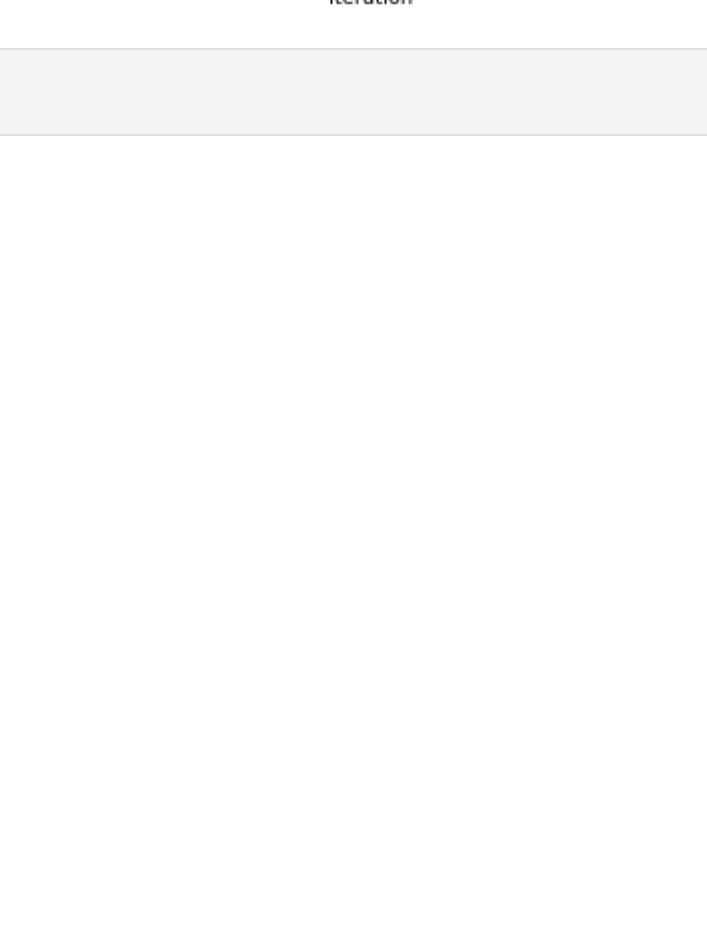
----- Q Learning (alpha=0.25, gamma=0.9, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.25, \gamma=0.9$



----- Q Learning (alpha=0.25, gamma=0.95, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.25, \gamma=0.95$



----- Q Learning (alpha=0.25, gamma=0.99, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.25, \gamma=0.99$



----- Q Learning (alpha=0.5, gamma=0.9, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.5, \gamma=0.9$



----- Q Learning (alpha=0.5, gamma=0.95, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.5, \gamma=0.95$



----- Q Learning (alpha=0.5, gamma=0.99, strategy: exponential) -----:

Average rate of success of the learned policy (Q Learning). $\alpha=0.5, \gamma=0.99$



A simple Count-based Exploration by using exponential decay:

$$\frac{dN}{dt} = -\lambda N$$

The solution to this equation is:

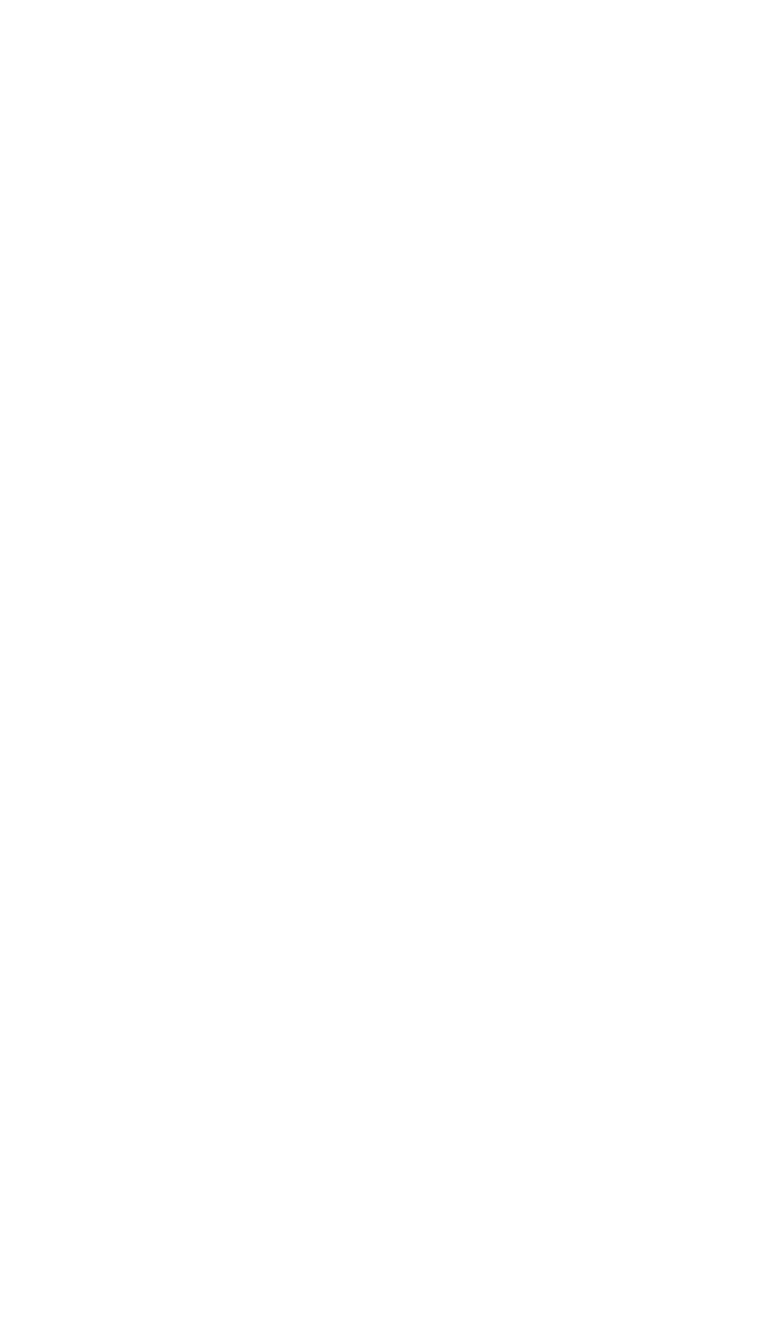
$$N(t) = N_0 e^{-\lambda t}$$

In this case N_0 is the initial exploration probability at $t = 0$ is 1 and the constant λ is the decay constant, $\lambda = 0.001$. Therefore, agent will absolutely explore at $t = 0$ and exponentially decay the exploration probability across time.

This exploration strategy does not work well in hard exploration problem but it is good enough for our environment.

```
In [18]: prob = []
init_prob = 1
lamd = 0.001
for t in range(5000):
    Nt = init_prob * np.exp(-lamd * t)
    prob.append(Nt)

plt.plot(prob)
plt.xlabel("Iteration")
plt.ylabel("Probability")
plt.title("Exploration Probability")
plt.show()
```



```
In [ ]:
```