

In [20]:

```
import sys
import time
import random
import numpy as np
import matplotlib.pyplot as plt
import sys
import gym
import numpy as np
import matplotlib.pyplot as plt
import sys
import gym
import numpy as np
import matplotlib.pyplot as plt
```

In [21]:

```
def simple_plot(
    arr,
    xlabel: str,
    ylabel: str,
    title: 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 show:
        plt.show()
    else:
        plt.savefig(save_path)
    if timeout:
        plt.pause(timeout)
    plt.draw()
    plt.waitforbuttonpress(timeout=5)
    plt.close()
    elif show:
        plt.show()
    plt.close()

# 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)
FFH       (H: hole, fall to your doom)
FFFG      (G: goal, where the frisbee is located)
print("Observation spaces ", env.observation_space)
print("Action space: ", env.action_space)

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

In [23]:

```
# Create MDP from the env as a reference
mdp = MDP(env)
```

In [24]:

```
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:.2f}")
    print(info, '\n')
    if done:
        break

Action sequence: ('Right', 'Right', 'Down', 'Down', 'Down', 'Right')
```

## 1.2

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:

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

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

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

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

$$\begin{aligned} v(s) &= E_{\pi} [G_t | S_t = s] \\ v(s) &= E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s \right] \\ &= E_{\pi} [R_t + \gamma \sum_{k=1}^{\infty} \gamma^{k-1} R_{t+k} | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma E_{\pi} \sum_{k=1}^{\infty} \gamma^{k-1} R_{t+k} | S_{t+1} = s'] \\ &= \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')] \end{aligned}$$

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

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

## 1.3

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  $\pi(s) = (s+1)$ .

In [25]:

```
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 i in range(trials):
        obs = self.env.reset()
        done = False
        while not done:
            obs, rew, done, info = self.env.step(act)
            reward += rew
            plt.show()
            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 rate: {success_rate}")
    return success_rate, mean_reward

# 3. Naive policy
policy = lambda s: (s + 1) % 4
naive_success_rates = []
for i in range(10):
    success_rate, mean_reward = 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: 0.013000000000000001
```

## 1.4

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

In [27]:

```
def LearnModel(self, n_samples: int = 10 ** 5) -> Tuple[np.ndarray, np.ndarray]:
    """
    Estimate transition probabilities  $p(s'|s, a, s')$  over  $n$  samples random samples.

    :param n_samples: Number of random samples to use.
    :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.nS, self.nA, self.nS)) # transition probability:  $S \times A \times S \rightarrow [0, 1]$ 
    R = np.zeros_like(P) # reward  $r(s, a, s')$ 

    obs = self.env.reset()
    for i in range(n_samples):
        # Random action
        act = self.action_space.sample()
        next_obs, rew, done, _ = self.env.step(act)
        P[obs, act, next_obs] += 1
        R[obs, act, next_obs] += rew
        obs = next_obs
        if done:
            obs = self.env.reset()
        # Normalize transition probabilities -> [0,1]
        pi = P.copy() # Don't modify P directly
        for a in range(self.nA):
            total_counts = np.sum(pi[s, a, :])
            if total_counts != 0:
                pi[s, a, :] /= total_counts
        # Avoid division by zero error
        R /= np.where(P == 0, 1, 1)

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

In [28]:

```
pi, r_hat = mdp.LearnModel(n_samples=10 ** 5)
mse = lambda x, y: np.mean((x - y) ** 2)

MSE_P, MSE_R = mse(mdp.P, pi), mse(mdp.R, r_hat)
print(f"Mean square error of P: {MSE_P}, \nMean square error of R: {MSE_R}")

Mean square error of P: 0.01957353902437787,
Mean square error of R: 0.0013020833333333333
```

## 1.5

Write a function `PolicyEval()` for evaluating a 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.

In [29]:

```
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 != None or self.R_hat != None:
        self.LearnModel()

    while True:
        delta = 0
        for s in range(self.nS):
            act = policy[s]
            V_s = self.P_hat[s, act, next_s] * (self.R_hat[s, act, next_s] + gamma * V[next_s])
            # Calculate delta
            delta = max(abs(V_s - V[s]), delta)
        # Update V
        V[:] = V_s
        if delta < theta:
            return V
        break
```

In [30]:

```
def PolicyIteration(
    self,
    max_iter: int = 50,
    gamma: float = 0.99,
    theta: float = 1e-6,
) -> 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), pi(s)
    V = np.zeros(self.nS)
    pi = np.zeros(self.nS, dtype=int) # since actions are integers
    success_rates = []
    mean_rewards = []

    print(f"Policy Iteration -----")
    for i in range(max_iter):
        P_hat, R_hat = self.LearnModel()
        # Policy Evaluation
        V = self.PolicyEval(V, P_hat, gamma, theta)
        # Policy Improvement
        pi = self.PolicyImprovement(V, gamma)
        P_hat, R_hat = self.LearnModel()
        success_rate, mean_reward = self.TestPolicy(pi, trials=100, render=False, verbose=True)
        success_rates.append(success_rate)
        mean_rewards.append(mean_reward)

        if np.all(pi["P_hat"] == P_hat):
            break
    return pi, V, success_rates, mean_rewards
```

Policy Iteration Converge less than 10 iterations.

In [31]:

```
# 5. Policy iteration
pi, V_pi, success_rates, mean_rewards = mdp.PolicyIteration(50, theta=sys.float_info.epsilon, exhaustive=False)
print(f"pi: {pi}")
print(f"V_pi: {V_pi}")

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()

----- Policy Iteration -----:
Iteration 1: Success rate: 0.0
Iteration 2: Success rate: 0.0
Iteration 3: Success rate: 0.1
Iteration 4: Success rate: 0.67
Iteration 5: Success rate: 0.71
Iteration 6: Success rate: 0.72
Iteration 7: Success rate: 0.78
Iteration 8: Success rate: 0.78

Policy is stable in 8 iterations
P: [0 3 3 0 0 0 0 3 1 0 0 0 2 1 0]
V_pi: [0.51276631 0.47393719 0.44811929 0.43499381 0.52815534 0.
0.3550636 0.
0.70497111 0.83785129 0.
0.70497111 0.83785129 0.]

Average rate of success of the learned policy (Policy Iteration)
Success rate
Iteration
```

Policy Iteration exhaust all 50 iterations.

In [32]:

```
# 6. Value Iteration
pi, V_pi, success_rates, mean_rewards = mdp.ValueIteration(50, theta=sys.float_info.epsilon, exhaustive=True)
print(f"pi: {pi}")
print(f"V_pi: {V_pi}")

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()

----- Value Iteration -----:
Iteration 1: Success rate: 0.0
Iteration 2: Success rate: 0.0
Iteration 3: Success rate: 0.13
Iteration 4: Success rate: 0.68
Iteration 5: Success rate: 0.71
Iteration 6: Success rate: 0.71
Iteration 7: Success rate: 0.72
Iteration 8: Success rate: 0.78
Iteration 9: Success rate: 0.78
Iteration 10: Success rate: 0.8
Iteration 11: Success rate: 0.74
Iteration 12: Success rate: 0.74
Iteration 13: Success rate: 0.82
Iteration 14: Success rate: 0.75
Iteration 15: Success rate: 0.72
Iteration 16: Success rate: 0.74
Iteration 17: Success rate: 0.7
Iteration 18: Success rate: 0.75
Iteration 19: Success rate: 0.75
Iteration 20: Success rate: 0.81
Iteration 21: Success rate: 0.72
Iteration 22: Success rate: 0.78
Iteration 23: Success rate: 0.68
Iteration 24: Success rate: 0.81
Iteration 25: Success rate: 0.73
Iteration 26: Success rate: 0.75
Iteration 27: Success rate: 0.69
Iteration 28: Success rate: 0.75
Iteration 29: Success rate: 0.75
Iteration 30: Success rate: 0.76
Iteration 31: Success rate: 0.74
Iteration 32: Success rate: 0.77
Iteration 33: Success rate: 0.69
Iteration 34: Success rate: 0.74
Iteration 35: Success rate: 0.79
Iteration 36: Success rate: 0.73
Iteration 37: Success rate: 0.79
Iteration 38: Success rate: 0.75
Iteration 39: Success rate: 0.73
Iteration 40: Success rate: 0.77
Iteration 41: Success rate: 0.78
Iteration 42: Success rate: 0.78
Iteration 43: Success rate: 0.6
Iteration 44: Success rate: 0.71
Iteration 45: Success rate: 0.72
Iteration 46: Success rate: 0.73
Iteration 47: Success rate: 0.74
Iteration 48: Success rate: 0.75
Iteration 49: Success rate: 0.77
Iteration 50: Success rate: 0.71

P: [0 3 3 0 0 0 0 3 1 0 0 0 2 1 0]
V_pi: [0.4548399 0.3938936 0.35943487 0.34182928 0.46805345 0.
0.31587021 0.
0.67864816 0.83239507 0.
0.67864816 0.83239507 0.]

Average rate of success of the learned policy (Value Iteration)
Success rate
Iteration
```

## 1.6

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.

In [33]:

```
def ValueIter(
    self,
    max_iter: int = 50,
    gamma: float = 0.99,
    theta: float = 1e-6,
) -> 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), pi(s)
    V = np.zeros(self.nS)
    pi = np.zeros(self.nS, dtype=int) # since actions are integers
    success_rates = []
    mean_rewards = []

    print(f"Value Iteration -----")
    for i in range(max_iter):
        P_hat, R_hat = self.LearnModel()
        # Value Iteration
        V = self.ValueEval(V, P_hat, gamma)
        # Policy Improvement
        pi = self.PolicyImprovement(V, gamma)
        P_hat, R_hat = self.LearnModel()
        success_rate, mean_reward = self.TestPolicy(pi, trials=100, render=False, verbose=True)
        success_rates.append(success_rate)
        mean_rewards.append(mean_reward)

        if np.all(pi["P_hat"] == P_hat):
            break
    return pi, V, success_rates, mean_rewards
```

In [34]:

```
# 6. Value Iteration
pi, V_pi, success_rates, mean_rewards = mdp.ValueIter(50, theta=sys.float_info.epsilon)
print(f"pi: {pi}")
print(f"V_pi: {V_pi}")

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()

----- Value Iteration -----:
Iteration 1: Success rate: 0.0
Iteration 2: Success rate: 0.0
Iteration 3: Success rate: 0.13
Iteration 4: Success rate: 0.28
Iteration 5: Success rate: 0.29
Iteration 6: Success rate: 0.46
Iteration 7: Success rate: 0.37
Iteration 8: Success rate: 0.39
Iteration 9: Success rate: 0.4
Iteration 10: Success rate: 0.6
Iteration 11: Success rate: 0.3
Iteration 12: Success rate: 0.37
Iteration 13: Success rate: 0.41
Iteration 14: Success rate: 0.48
Iteration 15: Success rate: 0.7
Iteration 16: Success rate: 0.71
Iteration 17: Success rate: 0.77
Iteration 18: Success rate: 0.73
Iteration 19: Success rate: 0.73
Iteration 20: Success rate: 0.72
Iteration 21: Success rate: 0.74
Iteration 22: Success rate: 0.74
Iteration 23: Success rate: 0.75
Iteration 24: Success rate: 0.79
Iteration 25: Success rate: 0.74
Iteration 26: Success rate: 0.72
Iteration 27: Success rate: 0.62
Iteration 28: Success rate: 0.73
Iteration 29: Success rate: 0.73
Iteration 30: Success rate: 0.73
Iteration 31: Success rate: 0.72
Iteration 32: Success rate: 0.77
Iteration 33: Success rate: 0.73
Iteration 34: Success rate: 0.73
Iteration 35: Success rate: 0.81
Iteration 36: Success rate: 0.72
Iteration 37: Success rate: 0.76
Iteration 38: Success rate: 0.82
Iteration 39: Success rate: 0.77
Iteration 40: Success rate: 0.74
Iteration 41: Success rate: 0.63
Iteration 42: Success rate: 0.76
Iteration 43: Success rate: 0.76
Iteration 44: Success rate: 0.73
Iteration 45: Success rate: 0.73
Iteration 46: Success rate: 0.73
Iteration 47: Success rate: 0.76
Iteration 48: Success rate: 0.77
Iteration 49: Success rate: 0.77
Iteration 50: Success rate: 0.77

P: [0 3 3 0 0 0 0 3 1 0 0 0 2 1 0]
V_pi: [0.4548399 0.3938936 0.35943487 0.34182928 0.46805345 0.
0.31587021 0.
0.67864816 0.83239507 0.
0.67864816 0.83239507 0.]

Average rate of success of the learned policy (Value Iteration)
Success rate
Iteration
```

(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  $\gamma$  increase, agent will focus more on the long-term reward as opposed to the immediate reward which is the case in FrozenLake-v0.

In [36]:

```
# (b)
for q in discount_factor:
    pi_QL, V_QL, success_rates_QL = mdp.QLearning(
        n_episodes=5000,
        gamma=q,
        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={q}"
        f"gamma={q}"
        f"gamma={q}"
    ),
    # save path=plot_dir / f"Q_L_{lr}_{q}_{g}.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
Success rate
Iteration
```

----- Q Learning (alpha=0.25, gamma=0.99, strategy: epsilon) -----:
Average rate of success of the learned policy (Q Learning), alpha=0.1, gamma=0.99
Success rate
Iteration

----- 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
Success rate
Iteration

(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  $\gamma$  increase, agent will focus more on the long-term reward as opposed to the immediate reward which is the case in FrozenLake-v0.

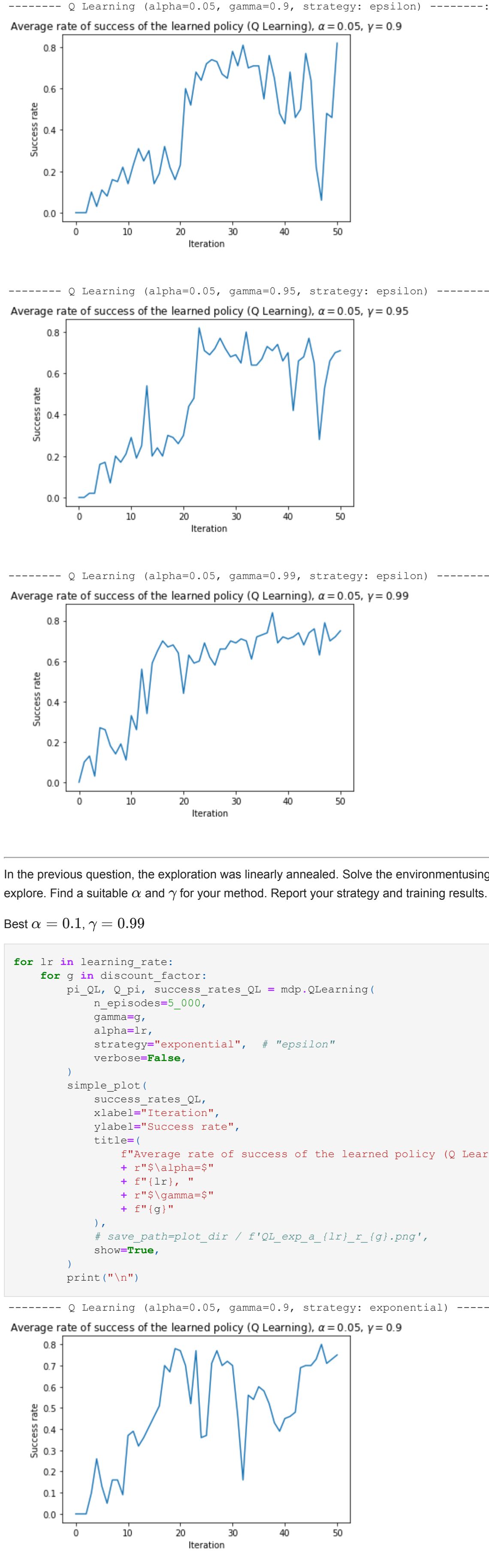
In [36]:

```
# (b)
for q in discount_factor:
    pi_QL, V_QL, success_rates_QL = mdp.QLearning(
        n_episodes=5000,
        gamma=q,
        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={q}"
        f"gamma={q}"
        f"gamma={q}"
    ),
    # save path=plot_dir / f"Q_L_{lr}_{q}_{g}.png",
    show=True,
    )
    print("\n")

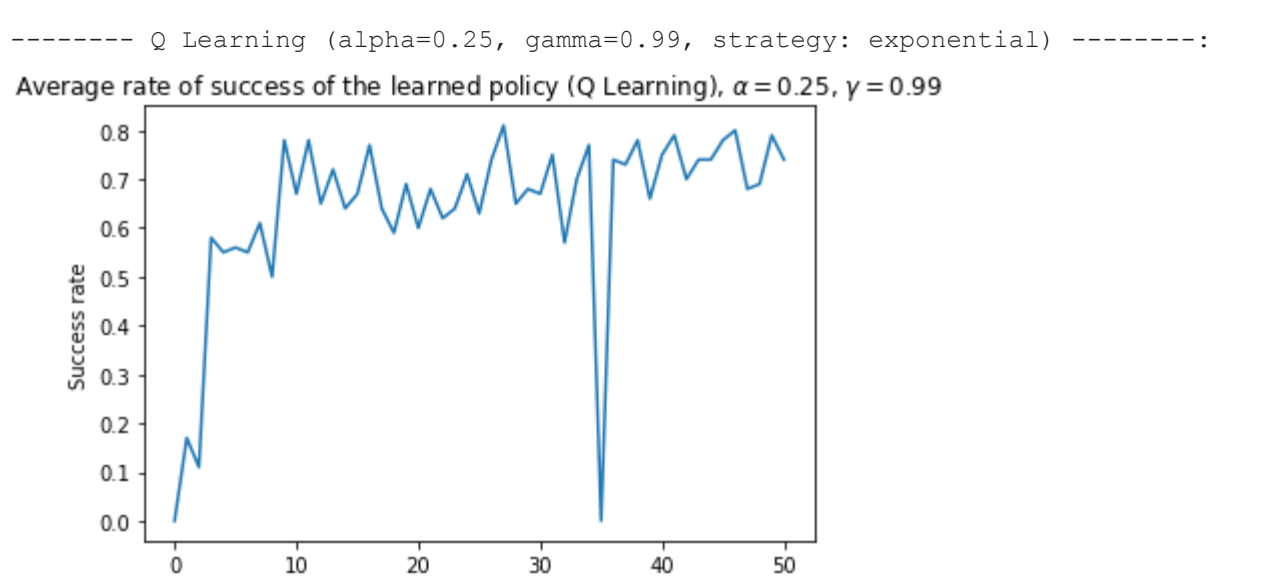
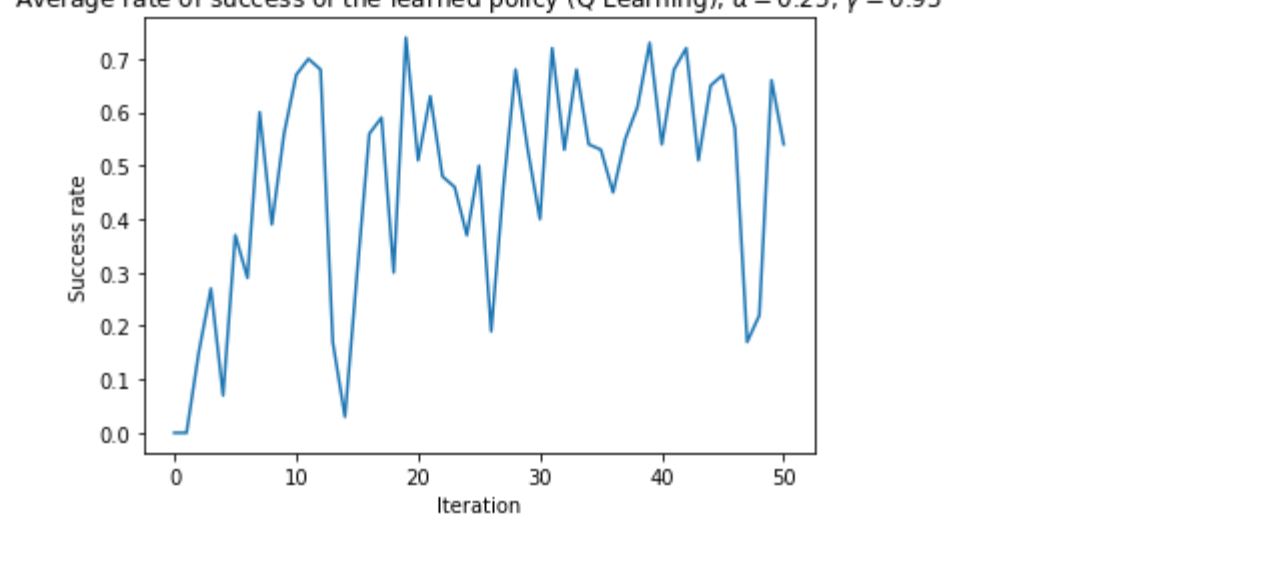
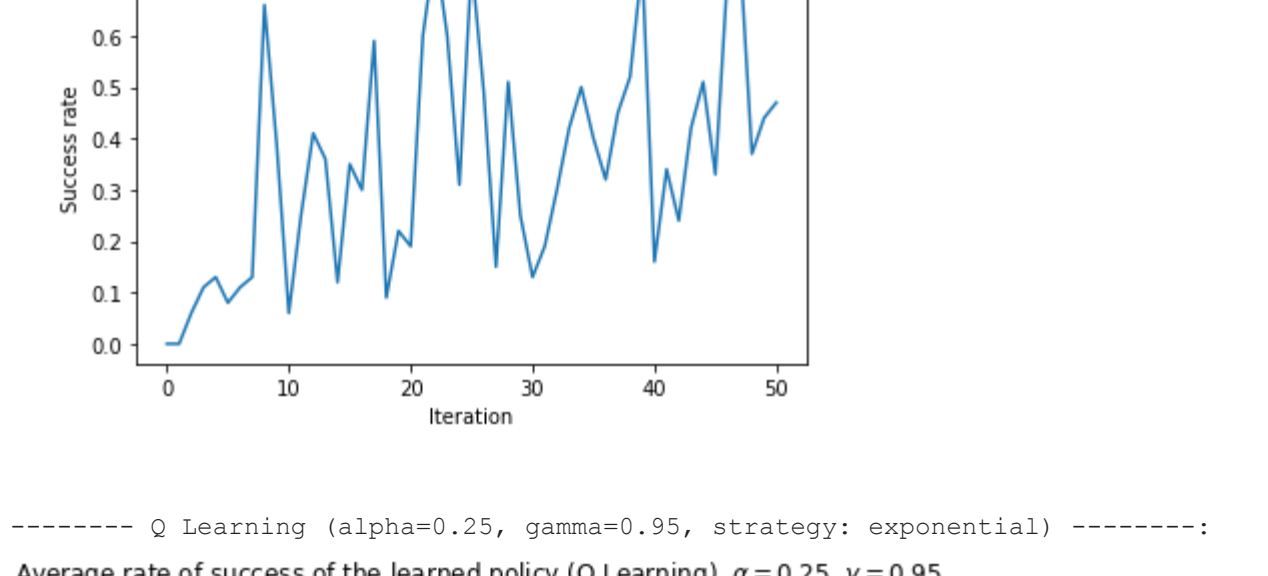
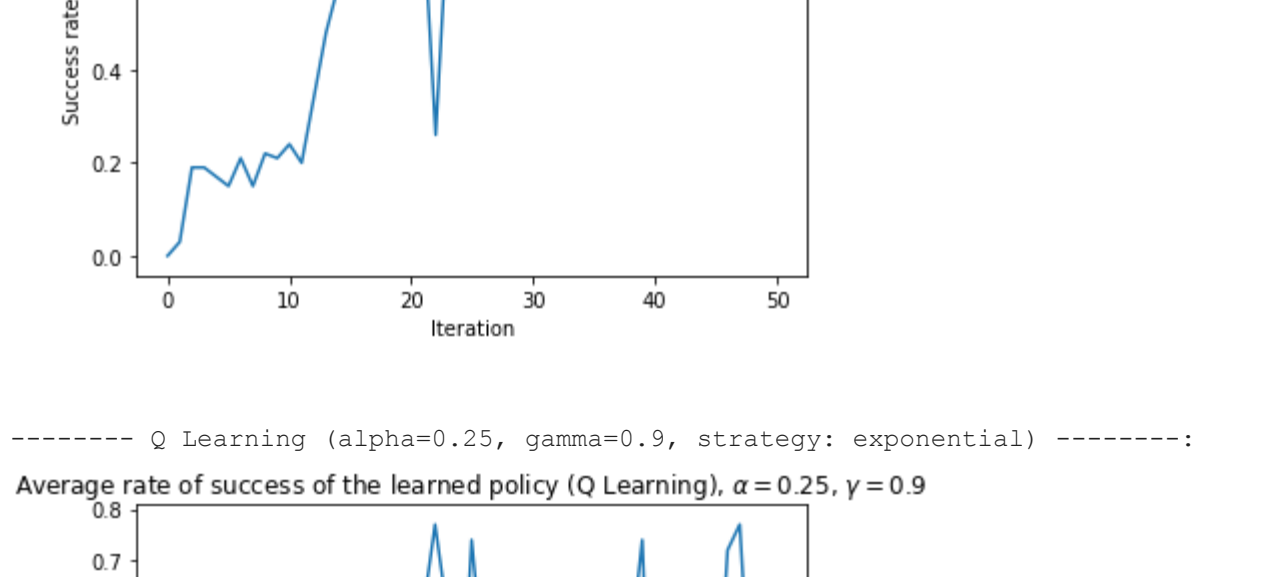
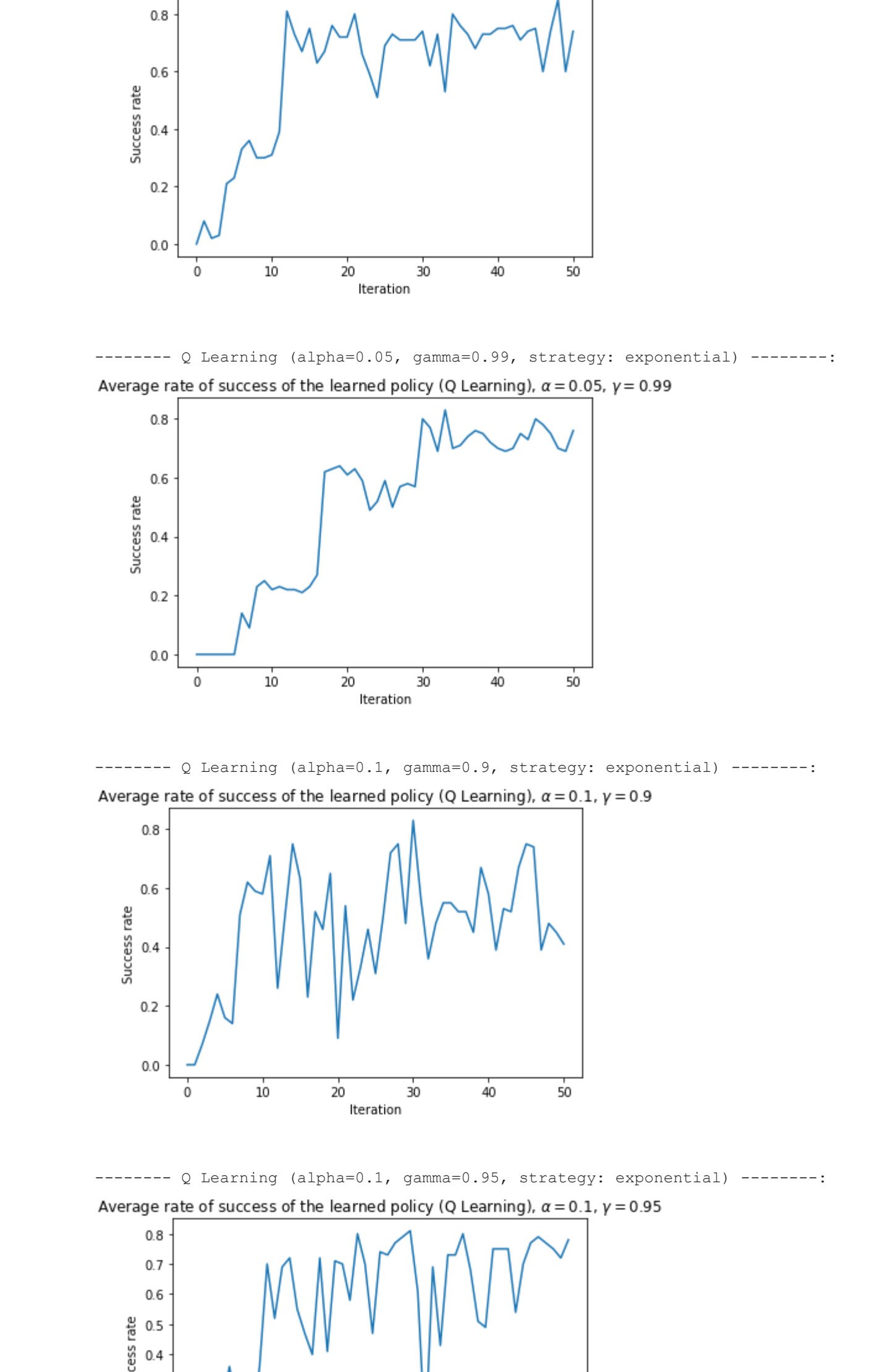
----- Q Learning (alpha=0.05, gamma=0.99, strategy: epsilon) -----:
Average rate of success of the learned policy (Q Learning), alpha=0.05, gamma=0.99
Success rate
Iteration
```





In the previous question, the exploration was linearily annealed. Solve the environment using Q-learning by proposing a different strategy to explore. Find a suitable  $\alpha$  and  $\gamma$  for your method. Report your strategy and training results.

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



A simple Count-based Exploration by using exponential decay:

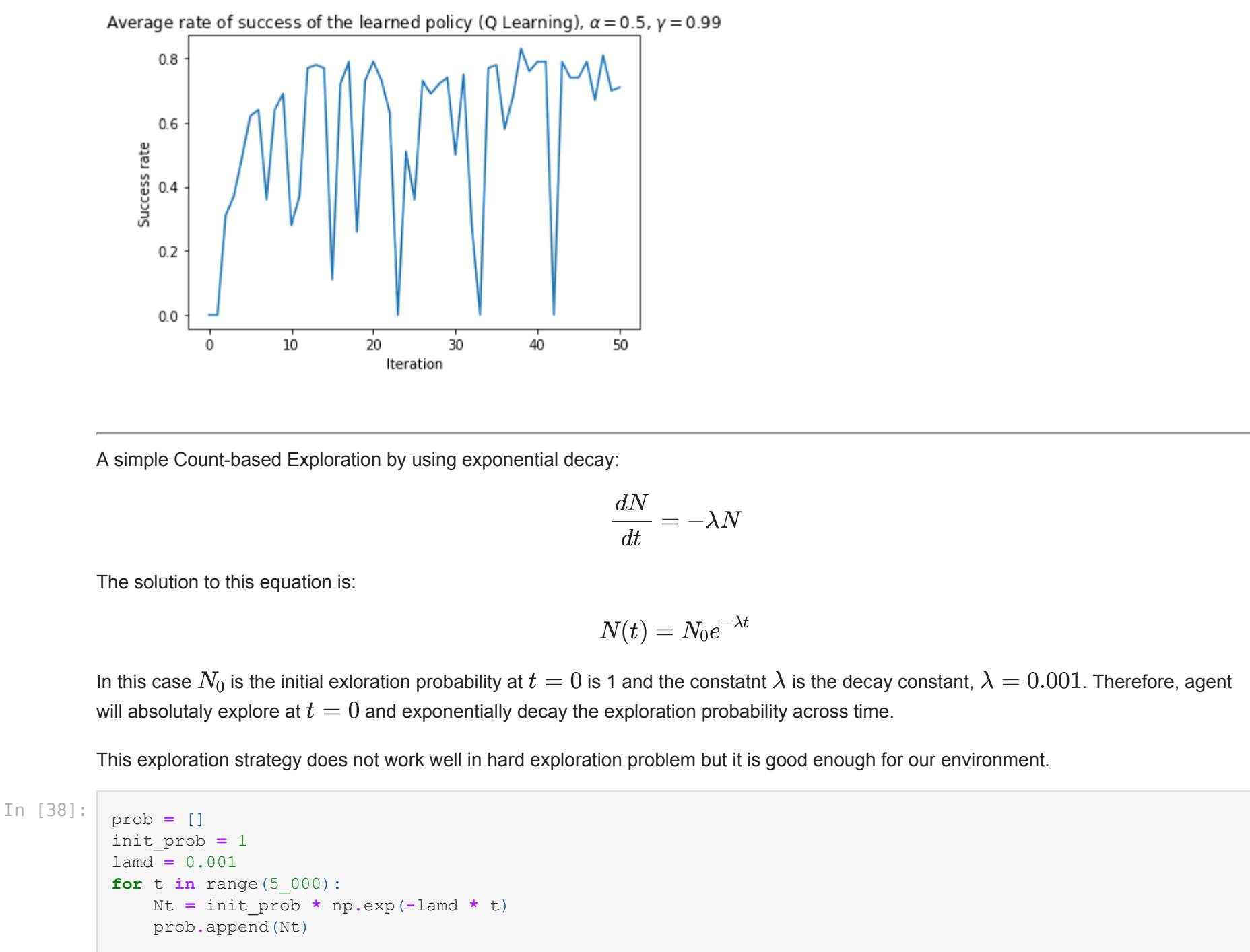
$$\frac{dN}{dt} = -\lambda N$$

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  $\lambda$  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 [ ]: