



SORBONNE UNIVERSITÉ

RAPPORT DE TME

Reinforcement Learning

Zhuangzhuang YANG 28708351
XIN HE 28706290
Master 1 Semestre 2

Tuteur : M. Olivier SIGAUD

Février 2021

2 MDPs and Mazes

2.2 Play with different MDPs

Code question 1

Here is the code which can build a maze like maze in figure 1.

```
1 walls = [7, 8, 9, 10, 21, 27, 30, 31, 32, 33, 45, 46, 47]
2 height = 6
3 width = 9
4 m = build_maze(width, height, walls) # maze-like MDP definition
```

3 Dynamic Programming

Dynamic programming algorithms are used for planning and they require a full knowledge of the MDP from agent. They can find the best policy by computing a value function V or an action-value function Q over the state space or state-action space of the given MDP.

3.1 Value iteration

To calculate the values following policy π with Q function, we implemented the Bellman Optimality Equation for Q^* . And here is the missing piece of code of `value_iteration_q(mdp)`:

Code question 3

```
1 for x in range(mdp.nb_states): #for each state x and action u
2     for u in mdp.action_space.actions:
3         if x in mdp.terminal_states:
4             # TODO: fill this
5             q[x, :] = mdp.r[x, u]
6         else:
7             # TODO: fill this
8             summ = 0
9             for y in range(mdp.nb_states):
10                 summ = summ + mdp.P[x, u, y] * np.max(q[y, :])
11             q[x, u] = mdp.r[x, u] + mdp.gamma * summ
```

The maze with the calculated state values :

0.17	0.19	0.21	0.19	0.17		0.48	0.53	0.59
0.15		0.23	0.21	0.19		0.53	0.59	0.66
0.17		0.25	0.23	0.21		0.59	0.66	0.73
0.19		0.28				0.53		0.81
0.21		0.31	0.35	0.39	0.43	0.48		0.9
0.23	0.25	0.28	0.31	0.35	0.39	0.43		1.0

3.2 Policy iteration

There are two steps of the Policy iteration :

- 1.evaluate policy π : compute V or Q based on the policy π
- 2.improve policy π : compute a better policy based on V or Q

We will use the Bellman Expectation with deterministic policy to find the optimal policy.

Code question 4

The function get_policy_from_q(q) :

```

1 def get_policy_from_q(mdp, q):
2     # Outputs a policy given the state values
3     policy = np.zeros(mdp.nb_states) # initial state values are set to 0
4     for x in range(mdp.nb_states): # for each state x
5         # Compute the value of the state x for each action u of the MDP action
        space
6         q_temp = np.zeros((mdp.nb_states,mdp.action_space.size))
7         for u in mdp.action_space.actions:
8             if x not in mdp.terminal_states:
9                 # Process sum of the values of the neighbouring states
10                summ = 0
11                for y in range(mdp.nb_states):
12                    summ = summ + mdp.P[x, u, y] * q[y,u]
13                q_temp[x,u] = mdp.r[x, u] + mdp.gamma * summ
14            else: # if the state is final, then we only take the reward into
                account
15                q_temp[x,:] = mdp.r[x, u]
16                policy[x] = np.argmax(q_temp)
17    return policy

```

Code question 5

The function evaluate_one_step_q(mdp,q,policy) :

```

1 def evaluate_one_step_q(mdp, q, policy):
2     # Outputs the state value function after one step of policy evaluation
3     # Corresponds to one application of the Bellman Operator
4     q_temp = np.zeros((mdp.nb_states, mdp.action_space.size))
5     for x in range(mdp.nb_states): # for each state x
6         # Compute the value of the state x for each action u of the MDP action
        space
7         for u in mdp.action_space.actions:
8             if x not in mdp.terminal_states:
9                 # Process sum of the values of the neighbouring states
10                summ = 0
11                for y in range(mdp.nb_states):
12                    summ = summ + mdp.P[x, u, y] * q[y,policy[y]]
13                q_temp[x,u] = mdp.r[x, u] + mdp.gamma * summ
14            else: # if the state is final, then we only take the reward into
            account
15                q_temp[x,:] = mdp.r[x, u]
16
17                # q_new[x] = np.max(q_temp)
18    return q_temp

```

The function evaluate_q(mdp,policy) :

```

1 def evaluate_q(mdp, policy):
2     # Outputs the state value function of a policy
3     q = np.zeros((mdp.nb_states, mdp.action_space.size)) # initial state values
        are set to 0
4     stop = False
5     while not stop:
6         qold = q.copy()
7         q = evaluate_one_step_q(mdp, qold, policy)
8
9         # Test if convergence has been reached
10        if (np.linalg.norm(q - qold)) < 0.01:
11            stop = True
12    return q

```

Code question 6

The missing code of the function policy_iteration_q(mdp) :

```

1 def policy_iteration_q(mdp, render=True): # policy iteration over the q
        function
2     q = np.zeros((mdp.nb_states, mdp.action_space.size)) # initial action
        values are set to 0
3     q_list = []
4     policy = random_policy(mdp)
5
6     stop = False
7
8     if render:
9         mdp.new_render()
10
11    while not stop:
12        qold = q.copy()
13
14        if render:
15            mdp.render(q)
16
17        # Step 1 : Policy evaluation

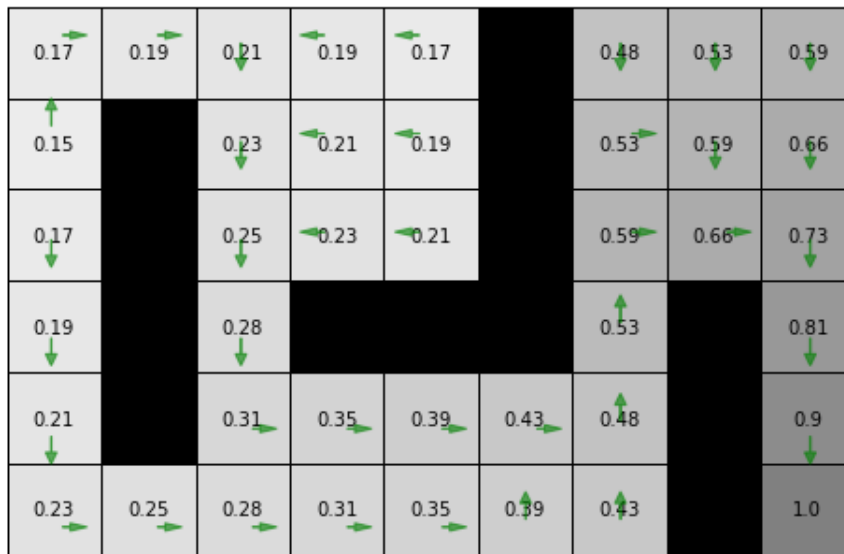
```

```

18     # TODO: fill this
19     q = evaluate_q(mdp,policy)
20
21     # Step 2 : Policy improvement
22     # TODO: fill this
23     policy = improve_policy_from_q(mdp, q, policy)
24     # Check convergence
25     if (np.linalg.norm(q - qold)) <= 0.01:
26         stop = True
27     q_list.append(np.linalg.norm(q))
28
29     # if render:
30     #     mdp.render(q, get_policy_from_q(q))
31     return q, q_list

```

And the final configuration :



Code question 7

The missing code of the fonction `policy_iteration_v(mdp)` :

```

1 def policy_iteration_v(mdp, render=True):
2     # policy iteration over the v function
3     v = np.zeros(mdp.nb_states) # initial state values are set to 0
4     v_list = []
5     policy = random_policy(mdp)
6
7     stop = False
8
9     if render:
10         mdp.new_render()
11
12     while not stop:
13         vold = v.copy()
14         # Step 1 : Policy Evaluation
15         v = evaluate_v(mdp,policy)
16
17         if render:
18             mdp.render(v)
19             mdp.plotter.render_pi(policy)
20

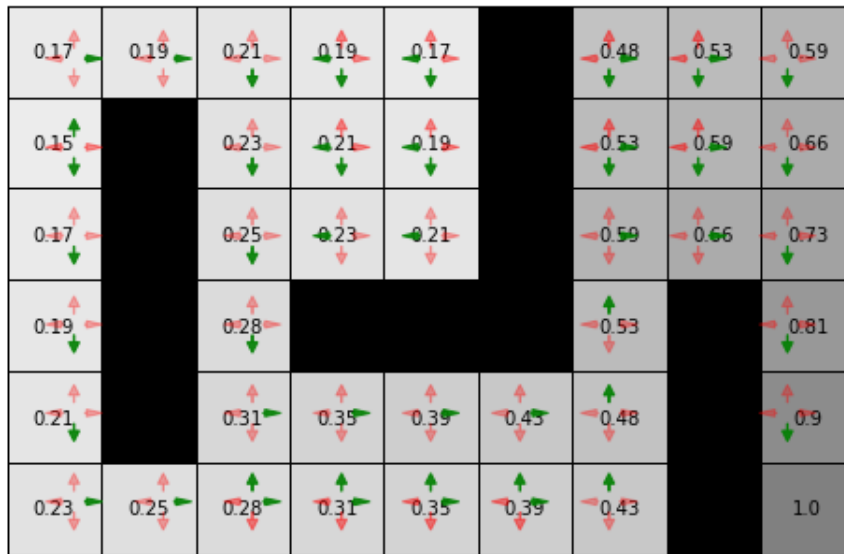
```

```

21     # Step 2 : Policy Improvement
22     # TODO: fill this
23     policy = improve_policy_from_v(mdp,v,policy)
24
25     # Check convergence
26     if (np.linalg.norm(v - vold)) < 0.01:
27         stop = True
28         v_list.append(np.linalg.norm(v))
29
30     if render:
31         mdp.render(v)
32         mdp.plotter.render_pi(policy)
33     return v, v_list

```

And the final configuration :



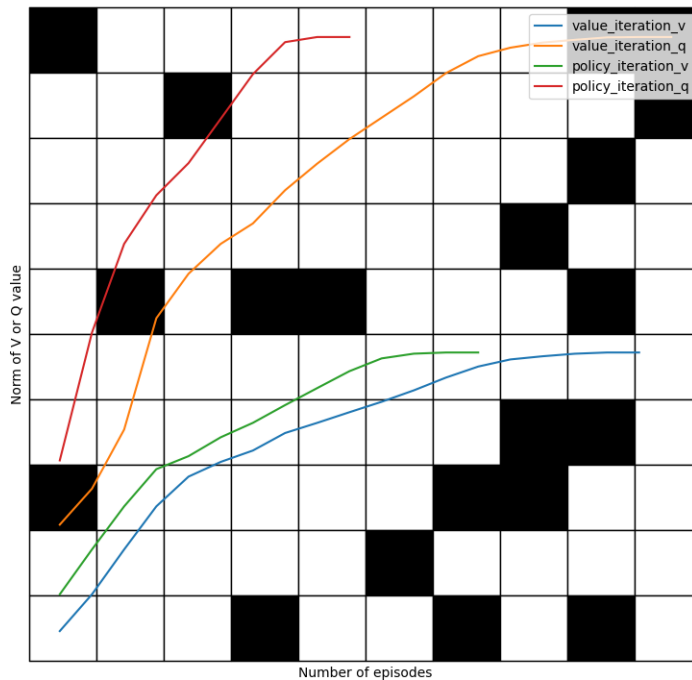
3.3 Comparisons

Study question 8

After add some code in dynamic programming functions, we have tested several different mazes, they have similar trends of results so that we select a represent one test :

method	nb_iterations	nb_elementary_updates	time_taken
value iteration V	480396	19	0s 380ms
value iteration Q	505680	20	2s 776ms
policy iteration V	13	155	1s 52ms
policy iteration Q	9	131	3s 53ms

And the image of convergence of those algorithms :



From the picture above, we can find the method `policy_iteration_q` has the fastest convergence, and the method `policy_iteration_v` has the second fastest convergence, but `value_iteration_v` and `value_iteration_q` have almost the same convergence. And we can find that `value_iteration_v` method take least time to get our final configuration, and it has less iterations and elementary than `value_iteration_q`, thus, we can drive the conclusion that `value_iteration_v` has better efficiency. For the two `policy_iteration` methods, `policy_iteration_v` runs faster but has a slower convergence, a more iterations of function and a more elementary updates. Thus `policy_iteration_v` is more efficient but not perfect. Generally, using a state value function is more efficient than a state-action value function.

4.1 TD learning

The TD-learning method is used for MDP which the agent doesn't know the transition and reward functions. It can compute the state value of the policy. And it's limitation : one cannot infer $\pi(s)$ from $v(s)$ without knowing T , one must know which action lead to best v' . And it's a On-policy method.

Code question 9

The corresponding code :

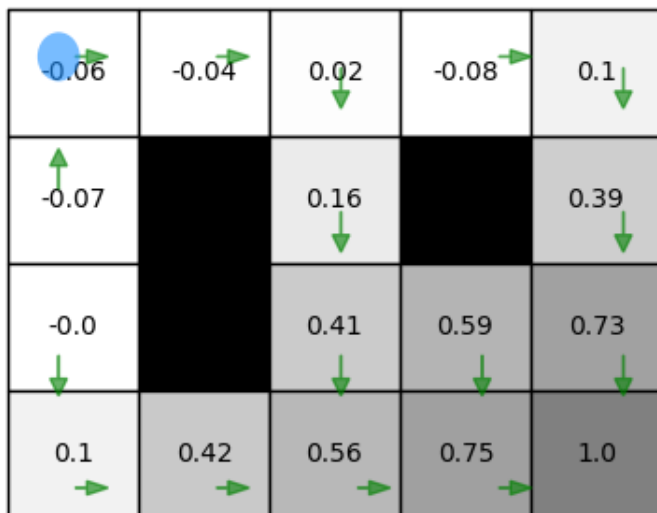
```
1 # Update the state value of x
2 if x in mdp.terminal_states:
```

```

3     v[x] = r
4     else:
5         delta = r + mdp.gamma * v[y] - v[x]
6         v[x] = v[x] + alpha * delta

```

And a picture of the final configuration :



4.2 Q-learning

Q-learning is a method that can solve the limitation of TD(0), an agent exploring an MDP and updating at each step a model of the state-action-value. No more need to know a_{t+1} unlike the method TD(0), and it's a off-policy method.

Code question 10

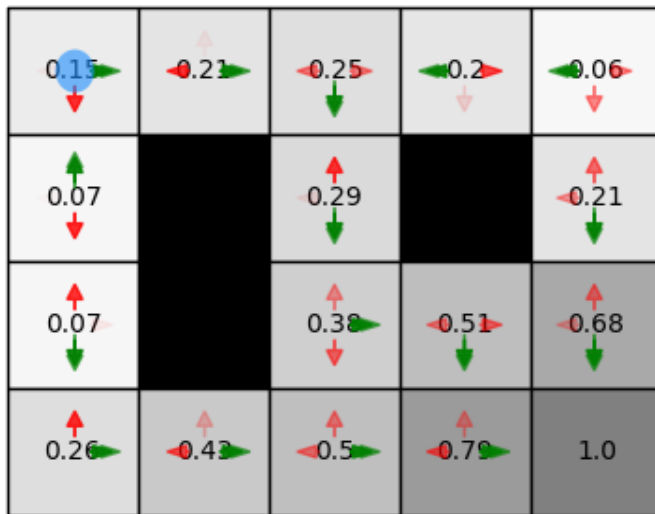
The corresponding code :

```

1     # Update the state-action value function with q-Learning
2     if x in mdp.terminal_states:
3         q[x, u] = r
4     else:
5         delta = r + mdp.gamma * np.max(q[y]) - q[x,u]
6         q[x, u] = q[x,u] + alpha*delta

```


And a picture of the final configuration :



Code question 11

The relevant piece of code :

```

1 def q_learning_eps(mdp, epsilon, nb_episodes=20, timeout=50, alpha=0.5, render=
  True):
2     # Initialize the state-action value function
3     # alpha is the learning rate
4     q = np.zeros((mdp.nb_states, mdp.action_space.size))
5     q_min = np.zeros((mdp.nb_states, mdp.action_space.size))
6     q_list = []
7     # Run learning cycle
8     mdp.timeout = timeout # episode length
9     if render:
10        mdp.new_render()
11    for _ in range(nb_episodes):
12        # Draw the first state of episode i using a uniform distribution over
13        all the states
14        x = mdp.reset(uniform=True)
15        done = mdp.done()
16        while not done:
17            if render:
18                # Show the agent in the maze
19                mdp.render(q, q.argmax(axis=1))
20            # Draw an action using a egreedy policy
21            u = egreedy(q, x, epsilon)
22
23            # Perform a step of the MDP
24            [y, r, done, _] = mdp.step(u)
25
26            # Update the state-action value function with q-Learning
27            if x in mdp.terminal_states:
28                q[x, u] = r
29            else:
30                delta = r + mdp.gamma * np.max(q[y]) - q[x,u]
31                q[x, u] = q[x,u] + alpha*delta
32            # Update the agent position
33            x = y
34        q_list.append(np.linalg.norm(np.maximum(q, q_min)))

```

```

34
35     if render:
36         # Show the final policy
37         mdp.current_state = 0
38         mdp.render(q, get_policy_from_q(q))
39     return q, q_list

```

4.3 SARSA

The difference between q_learning method is that Sarsa's update approach :

1. Sarsa is an update method of on-policy, and its action strategy and evaluation strategy are both ϵ -greedy strategies.
2. Sarsa is updated after first action, first execute the action through the ϵ -greedy strategy, and then update the value function according to the executed action.

Code question 12

```

1 #sarsa with softmax
2 def sarsa_soft(mdp, tau, nb_episodes=20, timeout=50, alpha=0.5, render=True):
3     # Initialize the state-action value function
4     # alpha is the learning rate
5     q = np.zeros((mdp.nb_states, mdp.action_space.size))
6     q_min = np.zeros((mdp.nb_states, mdp.action_space.size))
7     q_list = []
8
9     # Run learning cycle
10    mdp.timeout = timeout # episode length
11
12    if render:
13        mdp.new_render()
14
15    for i in range(nb_episodes):
16        print(i)
17        # Draw the first state of episode i using a uniform distribution over
18        # all the states
19        x = mdp.reset(uniform=True)
20        ux = 0
21        done = mdp.done()
22        while not done:
23            if render:
24                # Show the agent in the maze
25                mdp.render(q, q.argmax(axis=1))
26            # Draw an action using a soft-max policy
27            u = mdp.action_space.sample(prob_list=softmax(q, x, tau))
28            # Perform a step of the MDP
29            [y, r, done, _] = mdp.step(u)
30            # Update the state-action value function with q-Learning
31            if x in mdp.terminal_states:
32                q[x, u] = r
33            else:
34                uy = mdp.action_space.sample(prob_list=softmax(q, y, tau))
35                delta = r + mdp.gamma * q[y,uy] - q[x,u]
36                q[x, u] = q[x,u] + alpha*delta
37            # Update the agent position
38            x = y
39            q_list.append(np.linalg.norm(np.maximum(q, q_min)))

```

```

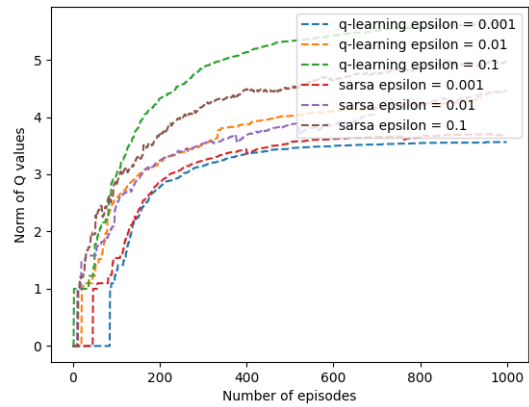
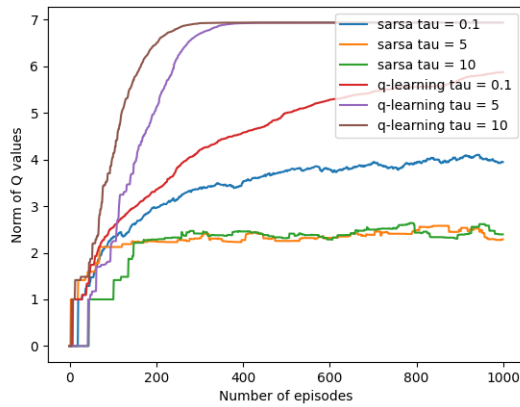
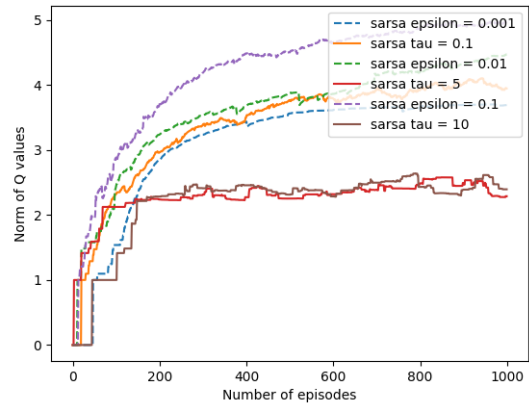
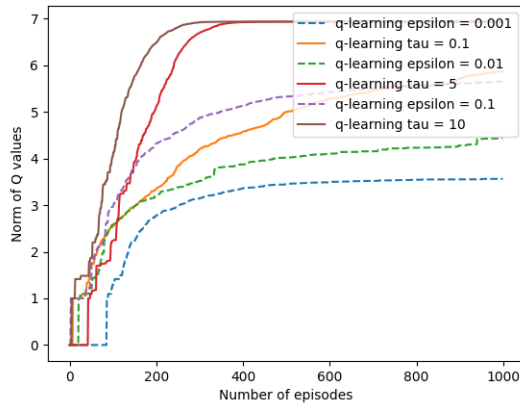
39     if render:
40         # Show the final policy
41         mdp.current_state = 0
42         mdp.render(q, get_policy_from_q(q))
43     return q, q_list
44
45 #sarsa with egreedy
46 def sarsa_eps(mdp, epsilon, nb_episodes=20, timeout=50, alpha=0.5, render=True):
47     # Initialize the state-action value function
48     # alpha is the learning rate
49     q = np.zeros((mdp.nb_states, mdp.action_space.size))
50     q_min = np.zeros((mdp.nb_states, mdp.action_space.size))
51     q_list = []
52     # Run learning cycle
53     mdp.timeout = timeout # episode length
54     if render:
55         mdp.new_render()
56     for i in range(nb_episodes):
57         print(i)
58         # Draw the first state of episode i using a uniform distribution over
59         # all the states
60         x = mdp.reset(uniform=True)
61         ux = 0
62         done = mdp.done()
63         while not done:
64             if render:
65                 # Show the agent in the maze
66                 mdp.render(q, q.argmax(axis=1))
67             # Draw an action using a egreedy policy
68             u = egreedy(q, x, epsilon)
69             # Perform a step of the MDP
70             [y, r, done, _] = mdp.step(u)
71             # Update the state-action value function with q-Learning
72             if x in mdp.terminal_states:
73                 q[x, u] = r
74             else:
75                 uy = egreedy(q, y, epsilon)
76                 delta = r + mdp.gamma * q[y,uy] - q[x,u]
77                 q[x, u] = q[x,u] + alpha*delta
78             # Update the agent position
79             x = y
80             q_list.append(np.linalg.norm(np.maximum(q, q_min)))
81         if render:
82             # Show the final policy
83             mdp.current_state = 0
84             mdp.render(q, get_policy_from_q(q))
85     return q, q_list

```

4.4 Comparisons and hyper-parameters

Study question 13

We used the function `plot_ql_sarsa` with various values of parameters `epsilon` and `tau` to get a image of curve below(the test come from the same maze) :

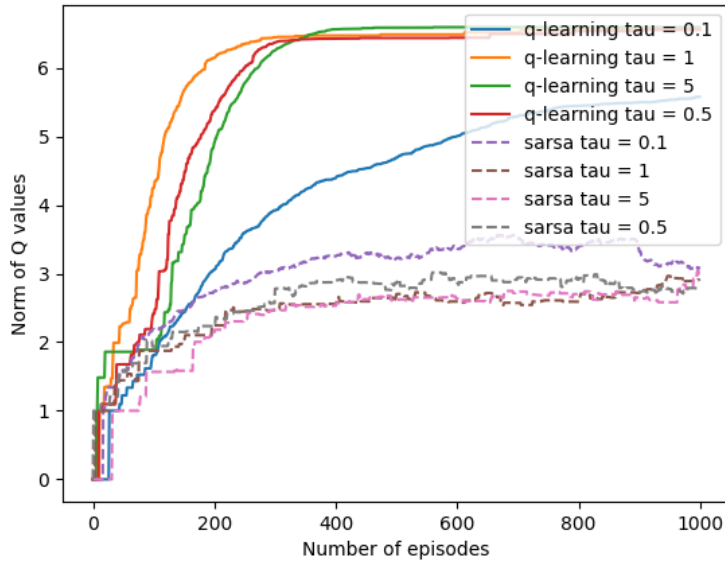


From the first image, we can see that the softmax method has faster convergence than ϵ -greedy, but if the $\epsilon=0.001$, it has the same performance as softmax method. From the second image, it's clear that softmax is better than ϵ -greedy because the red curve and brown curve have faster convergence, however the performance with $\epsilon=0.001$ is better than with $\tau=0.1$. And the third image, if $\tau=0.1$, we can find that sarsa has faster convergence, and it has the same situation when $\tau=5$ and $\tau=10$. However, in the final image, we can't drive a conclusion because the performance are the same when the parameters of epsilon are the same. Thus, we can drive several conclusions as follows :

1. The method sarsa has a faster convergence if we use softmax method.
2. The smaller epsilon is, the faster convergence will be.
3. Once the method q-learning has convergence, it's norm of q values won't change unlike the method sarsa.

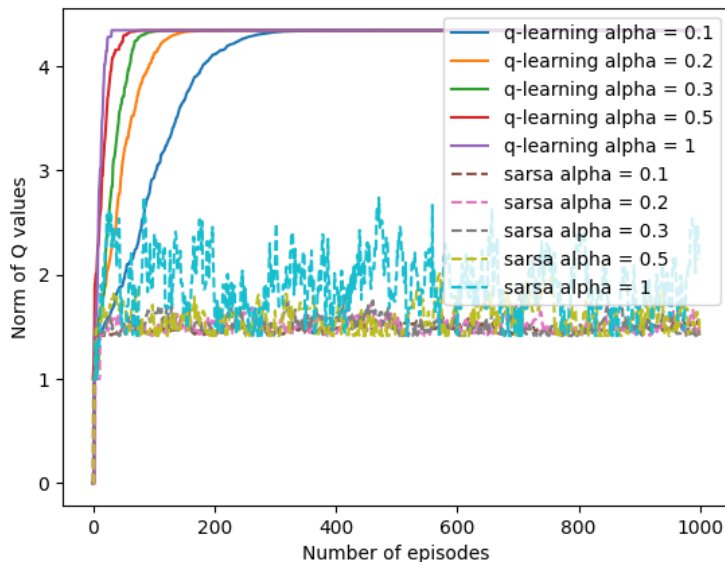
Study question 14

Here is the image of the relation of number of episode and Norm of Q value with parameter τ :



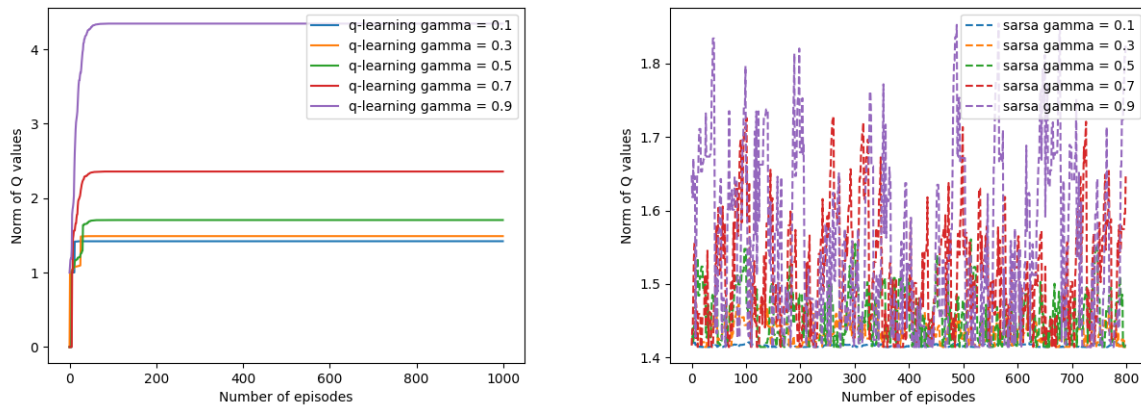
If τ is small, we observed the agent always move to the state where it can have the largest reward, so the agent can finish one episode rapidly, in this image, it represents the agent needs more episode to convergence like the blue curve. If the τ is large, the agent moved more randomly and passed by more state in one episode, so it represent the agent needs less episode to convergence, like the green curve in the image. However we observed that it runs faster if τ is smaller, so it is more efficient if the τ is smaller.

Here is the image of the relation of number of episode and Norm of Q value with parameter α :



If the value of α is smaller, the learning speed will be slower, we can see the blue curve has the slowest convergence with using q_learning method in this image. If the value of α is greater, the learning speed will be faster, we can see the purple one has the fastest convergence with using q_learning method in this image. Another influence : if the value of α is great, the state value changes evidently and the policy changes fast, sometimes we have already got the optimal policy, but it changed sooner, like those dotted curve in the image, their Q values changed fast, thus if the value of α is great, the policy is easier to oscillate.

Here is the image of the relation of number of episode and Norm of Q value with parameter γ :



The value of discount factor has an influence of MDP. To the q-learning method, we can see those curve take the same time to have convergence, the difference is that norm of q values is greater if discount factor is greater. And the oscillation is more strongly to the sarsa method if discount factor is greater. We have also taken an observation of the movement of agent, if the discount factor is closer to 0, the agent is more sensitive only to immediate reward, if the discount factor closer to 1, the future rewards are important as immediate reward.