

Internship report

Vo Van Nghia

Date

29 Sep, 2022

Table of contents

1	Institut de Recherche en Informatique de Toulouse (IRIT)	1
1.1	About the institut	1
1.2	Organization	1
2	System settings	1
2.1	Scheduling system	2
2.1.1	Parameters	2
2.1.2	Cost	2
2.1.3	Agent	2
2.1.4	Evolution	3
2.2	Load-balancing system	4
2.2.1	Parameters	4
2.2.2	Cost	4
2.2.3	Agent	4
2.2.4	Evolution	5
2.3	System representation	5
3	Reinforcement learning	6
3.1	Parameters	6
3.2	Bellman equation	6
3.3	Optimal policy for MDP	7
3.4	Efficient policy when having unobservable environment	8
4	Online Q-learning	8
4.1	About the algorithm	8
4.2	Implementation	9
4.2.1	From Python to C++	9
4.2.2	From n -dimensions to 2-dimensions	11
5	Offline Q-learning	12
5.1	About the algorithm	12
6	Value iteration	13
6.1	About the algorithm	13
6.2	Empirical transition probability	13
7	Result	14
7.1	No environment changing (MDP)	14
7.1.1	Scheduling	14
7.1.2	Load-balancing	15
7.2	Multiple environments (POMDP)	18
7.2.1	Scheduling	19
7.2.2	Load-balancing	23
8	Conclusion and extension	25

ACKNOWLEDGEMENT

I would like to thank my supervisors: Urtzi AYESTA, Maaïke VERLOOP and Matthieu JONCKHEERE for giving me an opportunity to work on this internship as well as your guidance and patience during the whole internship !

Abstract

Markov decision processes (MDPs) and their model free counterpart in reinforcement learning (RL) have known a large success in the last two decades. Although research in these two areas has been taking place for more than fifty years, the field gained momentum only recently following the advent of powerful hardware and algorithms with which suprahuman performance were obtained in games like Chess or Go. However, these impressive successes often rely on quite exceptional hardware possibilities and cannot be applied in many usual contexts, where, for instance, the volume of data available or the amount of computing power is more restricted. To define the next generation of more democratic and widely applicable algorithms, such methods still need to deal with very demanding exploration issues as soon as the state/action spaces are not small.

The internship aims to improve the efficiency of learning algorithms by leveraging the structure of the underlying problem and focus mainly on model free approach. The final report is divided into 8 parts. After a first introduction about the laboratory, we then present two situations where we study more carefully. Next, a short recap about MDPs and reinforcement learning, as well as discussions about three learning methods used in this report are presented. Finally, we compare the result of each method and have a final conclusion.

1 Institut de Recherche en Informatique de Toulouse (IRIT)

1.1 About the institut

The Institut de Recherche en Informatique de Toulouse (IRIT), created in 1990, is a Joint Research Unit (UMR 5505) of the Centre National de la Recherche Scientifique (CNRS), the Institut National Polytechnique de Toulouse (INP), the Université Paul Sabatier Toulouse3 (UT3), the Université Toulouse1 Capitole (UT1) and the Université de Toulouse Jean Jaurès (UT2J).

IRIT is one of the largest UMR at the national level, is one of the pillars of research in Occitanie with its 600 members, permanent and non-permanent, and about 100 external collaborators. The laboratory constitutes one of the structuring forces of the IT landscape and its applications in the digital world, both at regional and national level.

IRIT has focused its research on five major scientific issues and six strategic application areas.

- Health, Autonomy, Living, Well-being
- Smart City
- Aerospace and Transportation
- Social Media, Digital Social Ecosystems
- e-Education for learning and teaching
- Heritage and People Safety

As well as strategic action:

- Scientific Computing, Big Data and AI

My team is coming from team RMES (Réseaux, Mobiles, Embarqués, Sans fil, Satellites) and belongs to département ASR (Architecture, Systèmes, Réseaux).

1.2 Organization

The 24 research groups of the laboratory are dispatched in seven scientific departments:

- Dpt ASR : Architecture, Systems, Networks
- Dpt CISO : HPC, Simulation, Optimization
- Dpt FSL : Reliability of Systems and Software
- Dpt GD : Data Management
- Dpt ICI : Interaction, Collective Intelligence
- Dpt IA : Artificial Intelligence
- Dpt SI : Signals, Images

2 System settings

In this internship, we are interested in learning an efficient policy for some dynamic systems where the internal settings (transitions rate) of the system depend on some unobservable environments (or which environment the system is in), more generally, what is the optimal policy in the case of partial information. We could make one decision at a specific time (doing some things or even doing nothing) to make the system run more efficiently (faster or cheaper depending on each type of problem). However, we can only observe the state of the system (for example: the number of remaining tasks) and not the environments when making decisions. And therefore, we will explore Q-learning and derived techniques since they do

not need to know the internal settings of the systems. In this report, we focus ourselves on two following system models.

2.1 Scheduling system

2.1.1 Parameters

We have n classes of queues Q_1, \dots, Q_n , and L_1, \dots, L_n the maximum number of jobs (limit) on each class of queue. For each Q_i , its behavior is fully controlled by which environment it is in. An environment can be in one of the states m_1, \dots, m_m . And the environment of class Q_i is a random variable, denoted by M_i , which is in one of the states m_1, \dots, m_m . For each environment m_j , Q_i has their own holding cost $c_{i,j}$ (the cost of one unfinished job on the queue), arrival rate $\lambda_{i,j}$ (the rate of one more job arriving to the queue) and departure rate $\mu_{i,j}$ (the rate that a job departs the system (if served)). Furthermore, the environment M_i can change from m_j and m_k with the rate of $\xi_{i,j,k}$. A table summarizing all parameters is shown below.

Table 1: Scheduling system parameters

		m_1	...	m_m
Q_1	Holding cost	$c_{1,1}$...	$c_{1,m}$
	Arrival rate	$\lambda_{1,1}$...	$\lambda_{1,m}$
	Departure rate	$\mu_{1,1}$...	$\mu_{1,m}$
\vdots				
Q_n	Holding cost	$c_{n,1}$...	$c_{n,m}$
	Arrival rate	$\lambda_{n,1}$...	$\lambda_{n,m}$
	Departure rate	$\mu_{n,1}$...	$\mu_{n,m}$

Table 2: Matrix transition of the environment of Q_i

	m_1	...	m_m
m_1	$\xi_{i,1,1}$...	$\xi_{i,1,m}$
m_m	$\xi_{i,m,1}$...	$\xi_{i,m,m}$

The state of the system is represented by two vectors:

- The state $S = (X_1, \dots, X_n)$ where X_i is a random variable represents the current number of class Q_i jobs and is observable.
- The environment vector $E = (M_1, \dots, M_n)$ and this vector is not observable.

2.1.2 Cost

The cost of the system is a function of S and E . We choose a simple linear function for the cost:

$$f_1(S, E) = \sum_1^n c_{i, M_i} X_i$$

2.1.3 Agent

The agent will decide which queues should be activated. His goal is to minimize the cost of the whole system and the only information provided to him is the observable state S . Only jobs on activating queues can be processed and finished. In more generic problems, the agent is allowed to activate / deactivate multiple queues at the same time based on some conditions. However, in this internship, the agent can only activate one queue at a time.

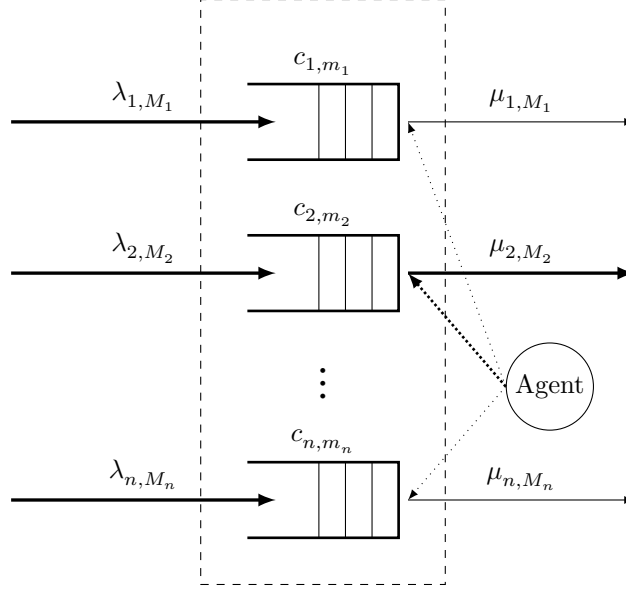


Figure 1: Visualization of the scheduling system when the agent activates queue 2, the bold lines represent the flow of jobs inside the system.

2.1.4 Evolution

In the continuous-time scale, each type of event (job arrival, job departure and environment changing) happens independently. That system is quite hard to program, therefore, we used a technique called “uniformization” to move from the continuous-time scale to discrete-time scale. In this time scale, at a given time, only one event can take place, regardless of their type.

In particular, given the system state is (S, E) , the agent decides to activate queue a . One of these $n + 1 + n(m - 1) + 1$ events can happen.

- The number of jobs in class Q_i increases by 1. Because there are n classes, we have n events of this kind with the rate $\lambda_{1,M_1}, \dots, \lambda_{n,M_n}$ respectively (if the number of job of class Q_i reaches an upper limit, we consider $\lambda_{i,M_i} = 0$).
- The number of jobs in class Q_a decreases by 1. Because we can only activate only one class, there is only one event of this type and its rate is μ_{a,M_a} .
- The environment of class Q_i changes to a different environment other than m_i and the rate changing to environment j is $\xi_{i,M_i,j}$. Because there are $m - 1$ possible changes for each class and there are n classes, the number of this kind of event is $n(m - 1)$.
- And a special dummy event where nothing changes.

A discrete probability distribution is used to express that. In order to satisfy the condition of a probability distribution, all the rates above are divided by a normalization constant C to make sure that their sum is not greater than 1. If that sum is smaller than 1, the special dummy event is used to fill the gap so that the final sum will be equal to 1.

The normalization constant has the form as follows, which is deduced from the above evolution of the system.

$$C = \sum_{i=1}^n \max_j \lambda_{i,j} + \max_{i,j} \mu_{i,j} + \sum_{i=1}^n \max_j \sum_{k=1, k \neq j}^m \xi_{i,j,k} + \epsilon$$

If $\epsilon > 0$, the probability of the dummy event will always be greater than 0. In this internship, as we do

not want the system to evolve too slowly, we choose $\epsilon := 0$.

After obtaining all the information above, we use that discrete probability distribution to obtain the next transition T of the system and denote S' and E' the next state of the system.

2.2 Load-balancing system

2.2.1 Parameters

For the load-balancing system, we only have one stream of job arrivals with rate λ and n queues. Same as above, The queues capacities (holding cost and departure rate) also depend on which environment they are in. Therefore, we have similar parameters to the one of the scheduling system. The only difference in this system is that instead of having multiple arrival rates that depend on the class of queue as well as the environment, we have only one global arrival rate λ . A table summarizing the parameters of this system is shown below.

Table 3: Load-balancing system parameters

		m_1	...	m_m
Q_1	Holding cost	$c_{1,1}$...	$c_{1,m}$
	Departure rate	$\mu_{1,1}$...	$\mu_{1,m}$
\vdots				
Q_n	Holding cost	$c_{n,1}$...	$c_{n,m}$
	Departure rate	$\mu_{n,1}$...	$\mu_{n,m}$
	Global arrival rate	λ		

2.2.2 Cost

We use the same cost function of scheduling system here.

2.2.3 Agent

The goal of the agent is the same: minimizing the cost of the whole system. He also receives the same information: the observable state S . The difference here is instead of having to deactivate some queues, all queues now run continuously. His mission is to choose which queue to send the new arriving job to. Only the chosen queue will have new jobs arriving, the other queues only have to process the remaining jobs.

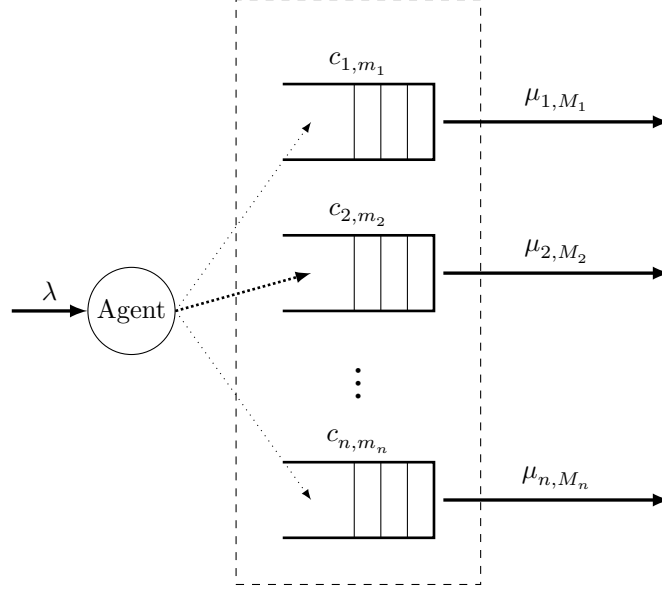


Figure 2: Visualization of the load-balancing system when the agent sends job to queue 2, the bold lines represent the flow of jobs inside the system.

2.2.4 Evolution

In this system, given its state is (S, E) , the agent decides to send jobs to queue a . One of these $1 + n + n(m - 1) + 1$ events can happen.

- The number of jobs in class Q_a increases by 1. Because we can only send jobs to one queue, there is only one event of this type and its rate is λ .
- The number of jobs in class Q_i decreases by 1. Because there are n classes, we have n events of this kind with the rate $\mu_{1,M_1}, \dots, \mu_{n,M_n}$ respectively (if the number of job of class Q_i reaches 0, we consider $\mu_{i,M_i} = 0$).
- The environment of class Q_i changes to a different environment other than m_i and the rate changing to environment j is $\xi_{i,M_i,j}$. Because there are $m - 1$ possible changes for each class and there are n classes, the number of this kind of event is $n(m - 1)$.
- And a special dummy event where nothing changes.

The form of the normalization constant also changes according to the evolution of the system.

$$C = \lambda + \sum_{i=1}^n \max_j \mu_{i,j} + \sum_{i=1}^n \max_j \sum_{k=1, k \neq j}^m \xi_{i,j,k} + \epsilon$$

2.3 System representation

In this internship, we limited ourselves on the cases where $n = 2$ for both systems. A system of that kind can be represented as a grid where each side represents the evolution of each queue. In this report, we use the vertical side for the first queue, and the horizontal one for the second. Each cell of that grid represents a specific observable state of the system.

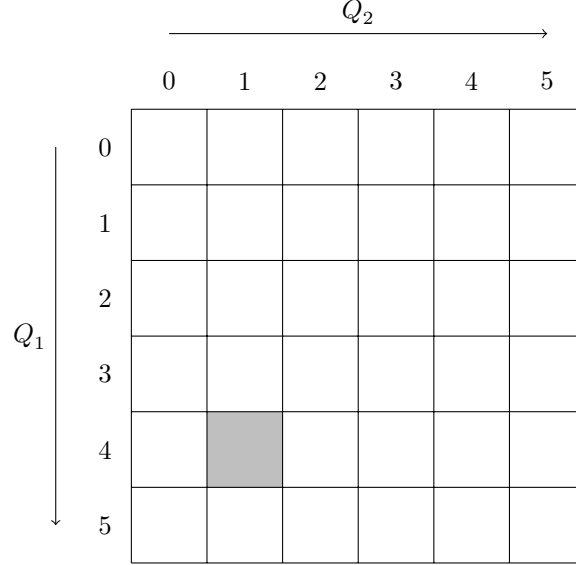


Figure 3: Representation of a system with 2 queues and the length limit of both queues are 5, the gray cell represents an observable state of the system $S = (4, 1)$.

3 Reinforcement learning

Reinforcement learning is used for training the agent to attain his goal. In this section, we present a short summary and introduce some related notions that could be useful for later.

3.1 Parameters

A simple reinforcement learning is modeled by a markov decision process (MDP) whose parameters include as follow:

- \mathcal{S} : the set of all states of the system.
- \mathcal{A} : the set of all actions, and A_S , the set of all actions available from state S .
- $P(S_{t+1} = S' | S_t = S, a_t = a)$: the probability that action a in state S at time t will lead to state S' at time $t + 1$.
- $R(S' | S, a)$: the immediate reward received after transitioning from state S to state S' , due to action a .

In addition, we have policy $\pi(a|S) = \mathbb{P}[A_t = a | S_t = S]$, a distribution over actions given states. In our problem, there is also an unobservable state inside our system, therefore, our problem is no longer an MDP, but instead a partial observable markov decision process (POMDP). We will now continue discussing the MDP and return to the POMDP in a later section.

3.2 Bellman equation

First, we want to try the solving method of a MDP.

The accumulate discounted future reward (return) from time t with the discount factor $0 \leq \gamma \leq 1$ is:

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots = r_{t+1} + \gamma G_{t+1}$$

From that, we have state-value function $V_\pi(S)$ is the expected return starting from state S , and then

following policy π :

$$V_\pi(S) = \mathbb{E}_\pi[G_t | S_t = S] \quad (1)$$

And action-value function $Q_\pi(S, a)$ is the expected return starting from state S , taking action a , and then following policy π :

$$Q_\pi(S, a) = \mathbb{E}_\pi[G_t | S_t = S, a_t = a] \quad (2)$$

From the fact that $\pi(a|S)$ is a distribution over \mathcal{A} given S , we have:

$$\begin{aligned} V_\pi(S) &= \mathbb{E}_\pi[G_t | S_t = S] \\ &= \sum_{a \in \mathcal{A}} \pi(a|S) \mathbb{E}_\pi[G_t | S_t = S, a_t = a] \\ &= \sum_{a \in \mathcal{A}} \pi(a|S) Q_\pi(S, a) \end{aligned} \quad (3)$$

Furthermore, we see that the expected future return of $Q_\pi(S, a)$ is the sum of the current reward $R(S)$ as well as the expected return of the next state regardless of action multiplied by the probability of moving to that state. Therefore, we have the equation below:

$$Q_\pi(S, a) = R(S) + \gamma \sum_{S' \in \mathcal{S}} P(S'|S, a) V_\pi(S') \quad (4)$$

Finally, the Bellman equations are obtained from (3) and (4):

$$\begin{aligned} V_\pi(S) &= \sum_{a \in \mathcal{A}} \pi(a|S) Q_\pi(S, a) \\ &= \sum_{a \in \mathcal{A}} \pi(a|S) (R(S) + \gamma \sum_{S' \in \mathcal{S}} P(S'|S, a) V_\pi(S')) \end{aligned} \quad (5)$$

$$\begin{aligned} &= R(S) + \gamma \sum_{a \in \mathcal{A}} \pi(a|S) \sum_{S' \in \mathcal{S}} P(S'|S, a) V_\pi(S') \\ Q_\pi(S, a) &= R(S) + \gamma \sum_{S' \in \mathcal{S}} P(S'|S, a) V_\pi(S') \\ &= R(S) + \gamma \sum_{S' \in \mathcal{S}} P(S'|S, a) \sum_{a' \in \mathcal{A}} \pi(a'|S') Q_\pi(S', a') \end{aligned} \quad (6)$$

3.3 Optimal policy for MDP

The core idea of this problem is to find an optimal policy π_* that maximizes the expected reward or $V(S)$. Mathematically, we have

$$\pi > \pi' \text{ if } V_\pi(S) > V_{\pi'}(S) \quad \forall S \in \mathcal{S}$$

And we want to find

$$\pi_* \text{ such that } \pi_* \geq \pi \quad \forall \pi$$

We can do that by finding the optimal state-value function $V_*(S) = \max_\pi V_\pi(S)$ or action-value function

$Q_*(S, a) = \max_{\pi} Q_{\pi}(S, a)$ and define the optimal policy as follow:

$$\pi_*(S) = \begin{cases} 1 & \text{if } a = \arg \max_{a \in \mathcal{A}} Q_*(S, a) \\ 0 & \text{otherwise} \end{cases}$$

Plugging all together into equation (5) and (6), we have:

$$V_*(S) = R(S) + \max_{a \in \mathcal{A}} \gamma \sum_{S' \in \mathcal{S}} P(S'|S, a) V_*(S') \quad (7)$$

$$Q_*(S, a) = R(S) + \gamma \sum_{S' \in \mathcal{S}} P(S'|S, a) \max_{a' \in \mathcal{A}} Q_*(S', a') \quad (8)$$

3.4 Efficient policy when having unobservable environment

Our problem is a POMDP, the parameters could be seen as follow:

- \mathcal{S} : the set of observable states is represented as the grid.
- \mathcal{E} : the set of unobservable environment states.
- \mathcal{A} : the action of activating / sending work to a queue.
- $P(S', E'|S, E, a)$: the transition probability that is clarified in the evolution section and depends on the both visible and hidden state (S, E).
- $R(S'|S, E, a)$: the immediate reward is $-C$ where C is the total holding cost of the system. Since this function depends only on the state (S, E) of the system, we could write $R(S, E)$ instead.

Since the goal of the internship is not about solving the POMDP, in the later sections, we will explore how learning algorithms of MDP work in the context of an unobservable environment.

4 Online Q-learning

4.1 About the algorithm

Online Q-learning is probably one of the most simple and popular algorithms for reinforcement learning problems. It is used to estimate $Q_*(S, a) \forall S \in \mathcal{S}$ and $a \in \mathcal{A}$ while interacting with the system. This algorithm only requires the real system or a simulator to interact with, and not the internal settings of that system, which makes it fit the context of this internship and more broadly real-life for a MDP, where these settings can change at any time. Beside its simplicity and close to real life context, the algorithm is also proven to converge almost surely when the number of visits to each state goes to ∞ in [WD92]. Visually, we can imagine that our agent will follow a grid trajectory of states and try to discover what is the action-value at that state with a specific action.

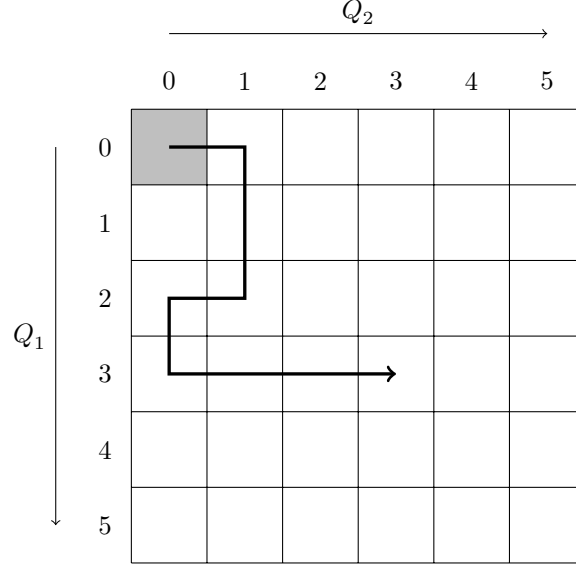


Figure 4: An example of the path which the agent might take. It starts from (0,0) and moves to (3,3). The next state could be either (3,4), (4,3), (3,2), (2,3) depending on the action it chooses and the transition probabilities.

In this report, we want to investigate the Q-learning algorithm when applied to our POMDP. We used the following algorithm:

Algorithm 1: Online Q-Learning

Input: $T > 0$ number iterations, γ discount factor, ϵ exploration factor, η learning rate

Output: Q_* , π_*

for $t \leftarrow 0$ **to** T **do**

$S, E \leftarrow \mathfrak{S}$ current state of the system;

$a \leftarrow \begin{cases} \text{one random possible action,} & \text{with the probability of } \epsilon \text{ (exploration)} \\ \arg \max_{a \in \mathcal{A}} Q_t(S, a), & \text{with the probability of } 1 - \epsilon \text{ (learning)} \end{cases}$

 /* $S' = S$ if $E' \neq E$ and vice versa */

$R, S', E' \leftarrow \mathfrak{S}(a, S, E);$

$Q_{t+1}(S, a) \leftarrow Q_t(S, a) + \eta(R + \gamma \max_{a' \in \mathcal{A}} Q_t(S', a') - Q_t(S, a));$

end

4.2 Implementation

4.2.1 From Python to C++

In this section, we document some difficulties while implementing the framework and how we overcome it. For illustration purposes, we will test the algorithm with a very simple problem of the **scheduling system** and linear cost function with no environment (or constant environment). The parameters of this system are shown in table 4.

It has been proven that if we have $c_1\mu_1 \geq c_2\mu_2$, the optimal policy for all cases is activating the first queue if possible, and vice versa in the case where $c_1\mu_1 \leq c_2\mu_2$. From table 4, we know that our optimal action here is activating the second queue.

With a naive Python implementation, the algorithm works correctly when the limits of both queues are small.

Table 4: Simple scheduling problem for testing

	Holding cost c_1	2
Q_1	Arrival rate λ_1	0.135
	Departure rate μ_1	0.3
	Holding cost c_2	3
Q_2	Arrival rate λ_2	0.135
	Departure rate μ_2	0.3

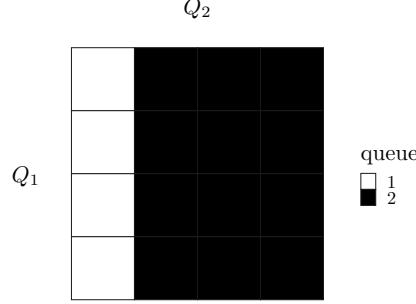


Figure 5: Result of online Q-learning using Python for 3x3 case. $\gamma = 0.999$, $\epsilon = 0.99$

First, we see that on the first column where $Q_2 = 0$, our agent always chooses to activate the first queue, because this is the only option possible, because there is no job on the second queue. Same thing holds for the first row where $Q_1 = 0$. For the other cells, we have the black color which means the agent chooses the second queue which is aligned with the analytical solution as above.

This implementation runs 10^6 iterations for approximately 2 minutes. This is quite good for systems with small limits. However, if we increase the limit to 10, with the same number of iterations gives an imperfect result.

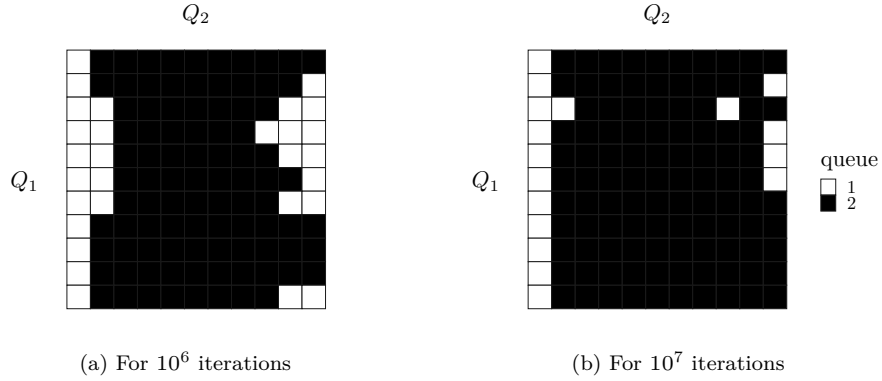


Figure 6: Result of online Q-learning using Python for 10x10 case. $\gamma = 0.99$, $\epsilon = 0.999$

In the figure above, we can see that even after 20 minutes with 10^7 iterations, the algorithm still does not converge completely yet. To have a clear view about the reason why, we have figure 7.

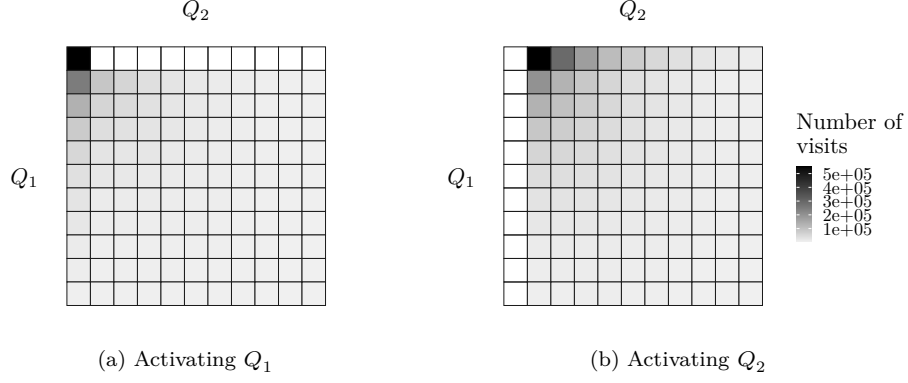


Figure 7: Number of visits for the system above. $\gamma = 0.99$, $\epsilon = 0.999$

We can see that our agent is stuck at $(0, 0)$, only the area around the origin has a high number of visits. The rest have relatively low number of visits and therefore, do not converge because of the nonfulfillment of the condition in [WD92].

To speed up the algorithm, we decided to switch to C++, which is famous for its speed and its mature support for high-speed computation with libraries like Eigen and then expose a Python interface for more ease of use. This is also the same approach for many scientific libraries in Python such as Numpy or Tensorflow. After switching to C++ and some improvements later, the algorithm now can run 10^9 iterations in only 1 minute, a 2×10^3 speed up.

4.2.2 From n -dimensions to 2-dimensions

With C++, the simulator is already running faster. However, there is still room for improvement. Now, the probability $P(S', E' | S, E, a)$ is calculated on-the-fly when we meet that state (S, E) while running the algorithm. We could make it faster by calculating the matrix transition for all the states beforehand. However, the state S and E of the system are each represented by a vector of n entries, if we keep using vector like that, we will need a tensor of $n + n + 1 = 2N + 1$ dimensions (first n dimensions for saving the observable state S , the next n dimensions for saving the unobservable state S and the last one for action). In addition, since there is no sparse version for high-dimensional tensor, it will take a lot of memory to save that tensor, more precisely:

$$(L_1 + 1) \times \cdots \times (L_n + 1) \times \underbrace{n^m}_{n \text{ queues and } m \text{ environments}} \times \underbrace{n}_{n \text{ actions}}$$

If we try to allocate a tensor of that size, we will get an “out of memory” error.

Furthermore, in other algorithms, we are required to iterate all over the set of states. If n , the number of queues is static, we could do something similar to algorithm 2. However our goal is to let n as a input for the program. So we can not use that algorithm there.

In the end, we decide to encode the whole state into an integer by extending the notion of base of number. Given state (S, E) , we introduce 3 functions:

$$\begin{aligned} f_S(S) &= f(X_1, \dots, X_n) = X_1 L_2 \dots L_n + X_2 L_3 \dots L_n + \cdots + X_{n-1} L_n + X_n \\ f_E(E) &= f(M_1, \dots, M_n) = M_1 m^{n-1} + M_2 m^{n-2} + \cdots + M_{n-1} m + M_n \\ f_{S,E(S,E)} &= f(X_1, \dots, X_n, M_1, \dots, M_n) = m^n f_S(S) + f_E(E) \end{aligned}$$

Algorithm 2: Iteration over the set of visible states \mathcal{S} when $n = 3$

```

for  $i \leftarrow 0$  to  $L_1$  do
  for  $j \leftarrow 0$  to  $L_2$  do
    for  $k \leftarrow 0$  to  $L_3$  do
       $S \leftarrow (i, j, k);$ 
      ...;
    end
  end
end

```

All three functions are one-to-one mapping, so we can replace the state with the result of these functions without losing anything. Because of doing so, we have achieved several benefits:

- We can calculate the transition matrix right now and use a sparse matrix to store it for efficiency in both storage and performance (since the transition matrix will be full of 0).
- Turn a problem with unknown dimension to a 2-dimensional one which makes accessing and looping through the states easier.
- Overall speed improvement since we only need to deal with integers instead of a vector.

5 Offline Q-learning

5.1 About the algorithm

In this version of Q-learning, we do not explore or follow the trajectory as in the online version, but we will force the calculation of the Q-value for all pairs (S, a) . This method is faster than online Q-learning, however, it requires a reliable simulator of the system to be able to work. The algorithm is described below:

Algorithm 3: Offline Q-Learning

Input: $T > 0$ number iterations, γ discount factor, η learning rate

Output: Q_*, π_*

\mathfrak{S} simulator of the system;

E initial environment state of the system;

```

for  $t \leftarrow 0$  to  $T$  do
   $-, -, E' \leftarrow \mathfrak{S}(-, -, E);$ 
  if  $E \neq E'$  then
     $E \leftarrow E';$ 
  end
  for  $S \in \mathcal{S}$  do
    for  $a \in \mathcal{A}$  do
       $S', R, - \leftarrow \mathfrak{S}(a, S, E);$ 
       $Q_{t+1}(S, a) \leftarrow Q_t(S, a) + \eta(R + \gamma \max_{a'} Q_t(S', a') - Q_t(S, a));$ 
    end
  end
end

```

6 Value iteration

6.1 About the algorithm

This algorithm is directly based on equation (7) and used to estimate the value of $V_*(S)$. After obtaining an estimate of $V_*(S)$, the policy is deduced as $\pi_*(S) = \arg \max_a R(S) + \gamma \sum_{S' \in \mathcal{S}} P(S'|S, a) V_*(S')$. The algorithm is written as follow:

Algorithm 4: Value iteration

Input: γ discount factor

Output: V_*, π_*

$\Delta \leftarrow \infty;$

while $\Delta > \epsilon$ **do**

for $S \in \mathcal{S}$ **do**

$V_{t+1}(S) \leftarrow \max_{a \in \mathcal{A}} R(S) + \gamma \sum_{S' \in \mathcal{S}} P(S'|S, a) V_*(S');$

end

$\Delta \leftarrow \|V_{t+1} - V_t\|_\infty;$

end

We can see above that the algorithm makes use of $P(S'|S, a)$, because of this, this algorithm converges a lot faster than both Q-learning methods but we have to know the transition probability, so this method is used only as a baseline to compare with the result of Q-learning.

If we have only one constant environment E_c (or no environment changing), we could rewrite $P(S', E_c|S, E_c, a) = P(S'|S, a)$ and value iteration work correctly. However, as we have unobservable environment, algorithm 4 can not work without modification because we have $P(S', E'|S, E, a)$ but we want to estimate $V_*(S)$ and not $V_*(S, E)$. In the next section, we discuss about how to get an empirical estimate of $P(S'|S, a)$ based on $P(S', E'|S, E, a)$.

6.2 Empirical transition probability

Our idea is to calculate the probability of the system stay in a particular state $E \in \mathcal{E}$ and then multiply it with the transition probability to get the average probability. First, we have some equations:

$$P(S'|S, E, a) = \begin{cases} P(S', E|S, E, a) & \text{if } S' \neq S \\ 1 - \sum_{S' \neq S} P(S', E|S, E, a) = \sum_{E'} P(S, E'|S, E, a) & \text{otherwise} \end{cases} \quad (9)$$

$$P(E'|S, E, a) = \begin{cases} P(S, E'|S, E, a) & \text{if } E' \neq E \\ 1 - \sum_{E' \neq E} P(S, E'|S, E, a) & \text{otherwise} \end{cases} \quad (10)$$

The logic behind equation (9) and (10) is the fact that only one event can happen at a time.

$$P(S, E'|S, E, a) = \begin{cases} \xi_{i, M_i, M'_i} & \text{if } E' \text{ and } E \text{ differ by only one entry at index } i \\ 0 & \text{else} \end{cases} \quad (11)$$

By rewriting $P(S, E'|S, E, a)$ as in equation (11), we see that $P(E'|S, E, a)$ does not depend on S, a and therefore $P(E'|S, E, a) = P(E'|E)$. From that, we could calculate the steady state by an iterative method using the equation below:

$$\pi_{t+1}(E') = \sum_{E \in \mathcal{E}} P(E'|E)\pi_t(E) \quad (12)$$

After $\pi(E)$ converges to $\tilde{\pi}(E)$ and using (9), we have the final empirical value we want:

$$\tilde{P}(S'|S, a) = \sum_{E \in \mathcal{E}} \tilde{\pi}(E) P(S'|S, E, a) \quad (13)$$

The empirical value we get here will be used in 4 in the case of POMDP for value iteration method.

7 Result

7.1 No environment changing (MDP)

7.1.1 Scheduling

First we will try with the linear cost. This is the simplest case and since it is proven mathematically in [BVW85], it is only a test to show that all algorithms work correctly. The optimal policy here will be activating the queue whose $c_i\mu_i$ is bigger. We will try two tuples of parameters, the first one will be:

Table 5: Simple scheduling 1

	Holding cost c_1	3
Q_1	Arrival rate λ_1	0.13
	Departure rate μ_1	0.3
	Holding cost c_2	1
Q_2	Arrival rate λ_2	0.13
	Departure rate μ_2	0.3

The theoretical optimal action here is 1.

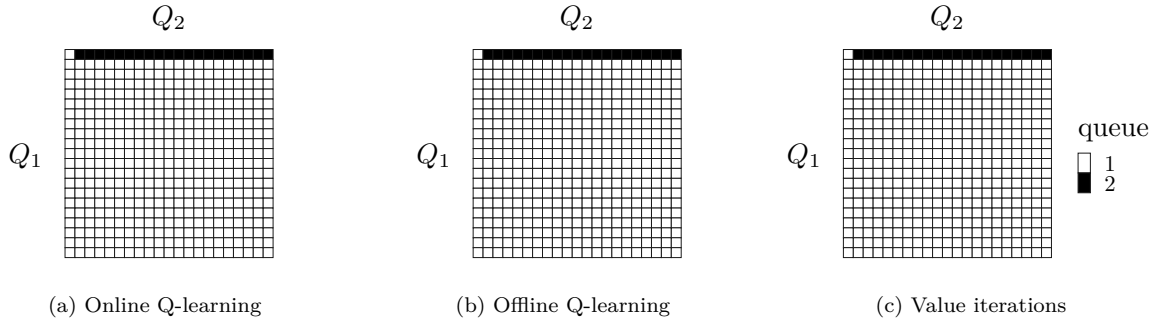
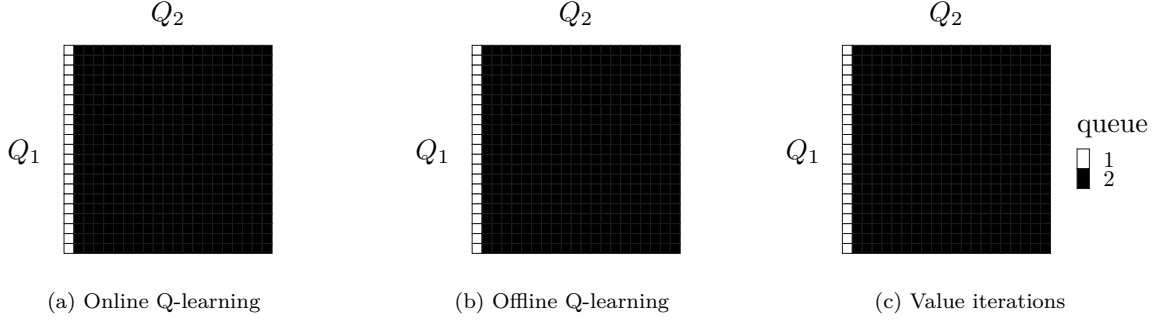


Figure 8: Result of table 5 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

And the second one is nearly the same, except the holding cost is reverse, and therefore the optimal action here is 2:

Table 6: Simple scheduling 2

Q_1	Holding cost c_1	1
	Arrival rate λ_1	0.13
	Departure rate μ_1	0.3
Q_2	Holding cost c_2	3
	Arrival rate λ_2	0.13
	Departure rate μ_2	0.3

Figure 9: Result of table 6 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

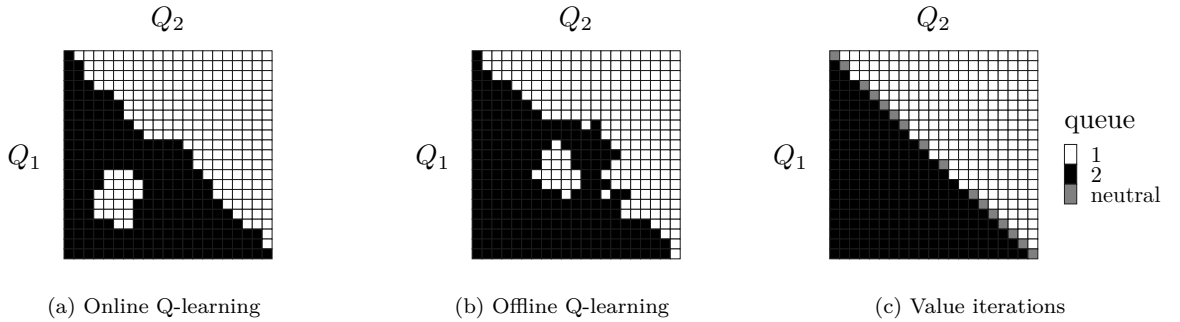
As we can see, all three algorithms work well for this simple case. We will move to more complex cases.

7.1.2 Load-balancing

For this system, We first try a system where two queues have the same parameters:

Table 7: Simple load-balancing 1

Q_1	Holding cost c_1	3
	Departure rate μ_1	0.3
Q_2	Holding cost c_2	3
	Departure rate μ_2	0.3
	Global arrival rate	0.5

Figure 10: Result of table 7 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

Intuitively, if two queues have the same parameters, the optimal action will be sending jobs to the queue that has fewer jobs, which is the result of the value iteration method (the right grid) of figure 10. There is also a gray color indicating that in those states, either action 1 (sending a new job to Q_1) or 2 (sending a new job to Q_2) leads to the same result, because they have the same number of jobs.

For two Q-learning methods, the online and offline version has been trained for 3×10^{10} and 5×10^7 iterations respectively and takes 30 minutes each. They give us a similar result and both are close to the optimal policy of value iteration.

However, since both methods depend on the evolution of the system, they are very sensitive to the parameters of systems. For example, if we increase the arrival rate to 0.55:

Table 8: Simple load-balancing 2

	Holding cost c_1	3
Q_1	Departure rate μ_1	0.3
	Holding cost c_2	3
Q_2	Departure rate μ_2	0.3
	Global arrival rate	0.55

Although the optimal policy should be the same as above, the same number of iterations give us a different result.

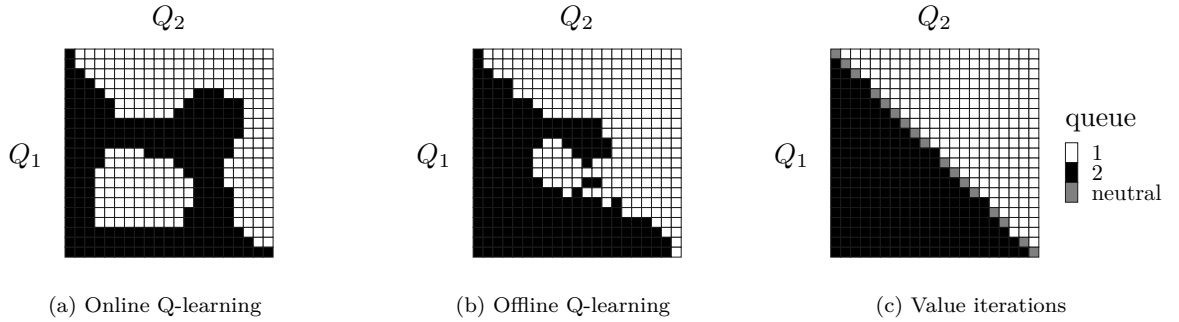


Figure 11: Result of table 8 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

The offline version is similar to the policy of value iteration method and has a smaller hole compared to the online version. The difference can be explained by the fact that in the offline version, each pair state-action is visited 5×10^7 times, so the total number of iterations will be $2 \times 41^2 \times 5 \times 10^7 \approx 1.6 \times 10^{11}$, which is $\gg 3 \times 10^{10}$, the total number of iterations of the online version.

Furthermore, even if we run the online version with that many iterations, it will still not guarantee that we will have a good approximation of the offline result. This is because we are following the trajectory (or the evolution of the system) and therefore it is impossible to visit all states equally. We can see the number of visits of each pair state-action below.

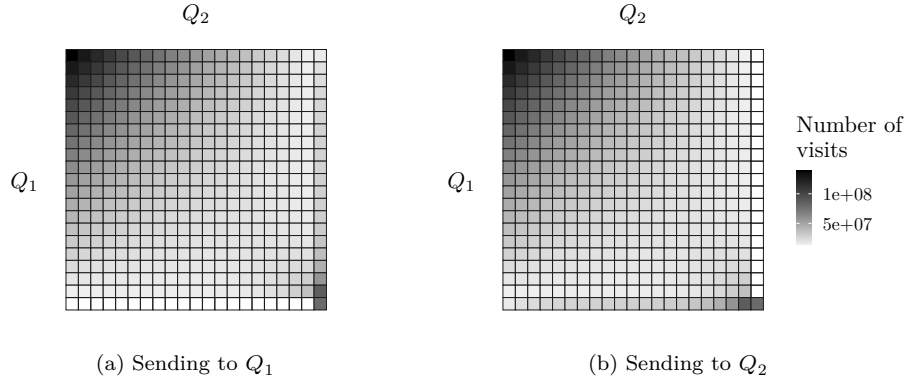


Figure 12: Number of visits of online Q-learning on table 8 for 20x20 case.

The number of visits is high around $(0,0)$ and $(40,40)$. Central region of the grid receives much less visits. Therefore, we see in figure 10 that the diagonal is beginning to form from these two points but not the region near the grid center.

We will finish this section with an unbalanced example. We fix the cost and make the departure rate of the first queue higher than the second one as follow:

Table 9: Simple load-balancing 3

	Holding cost c_1	3
Q_1	Departure rate μ_1	0.4
	Holding cost c_2	3
Q_2	Departure rate μ_2	0.3
	Global arrival rate	0.6

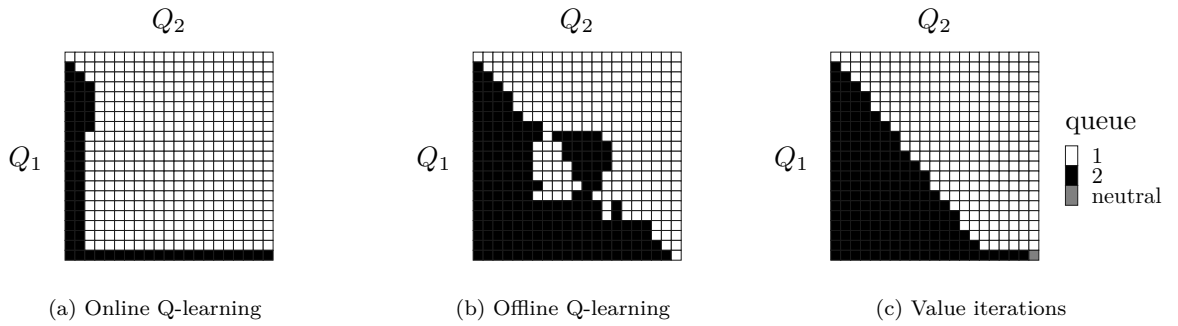


Figure 13: Result of table 9 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

Once again, we see the sensitivity of Q-learning with the system settings. Intuitively, since class 1 works faster than class 2 with the same cost, we will send more jobs to the first one as in the value iteration method. However, the learning capacity of both Q-learning methods is not sufficient to capture all the information of the system.

Furthermore, if we run the same setting but with 2×10^{11} iterations, the policy remains nearly the same although the $\|\cdot\|_2$ between two Q-values is approximately 3000. Since the difference is not small, we can see that the Q-value has not converged yet and we still need to run the simulation a lot longer.

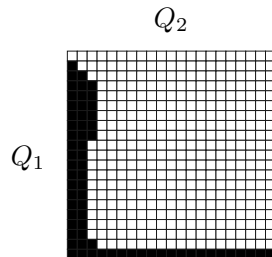
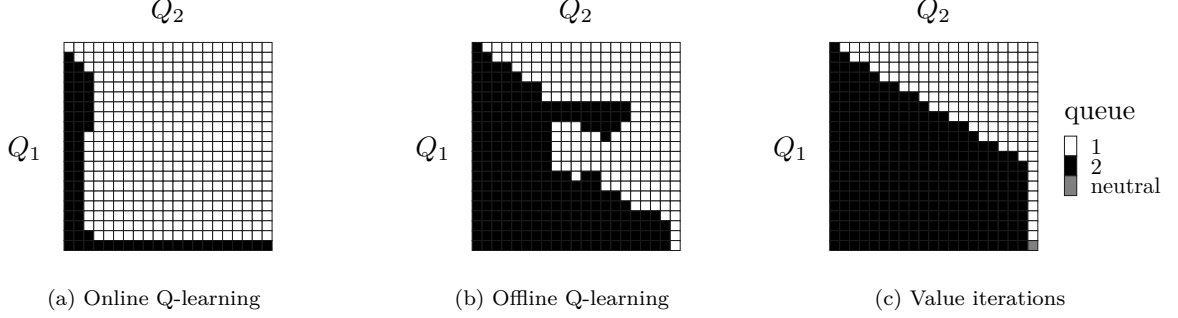


Figure 14: Result of table 9 for 20x20 case after 2×10^{11} iterations. $\gamma = 0.999$, $\epsilon = 0.999$.

And finally, we try a reversed setting where more jobs should be sent to queue 2 than queue 1.

Table 10: Simple load-balancing 4

Q_1	Holding cost c_1	3
	Departure rate μ_1	0.3
Q_2	Holding cost c_2	3
	Departure rate μ_2	0.5
	Global arrival rate	0.7

Figure 15: Result of table 9 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

After all the results above, we can see that generally, Q-learning performs badly on load-balancing systems even in the MDP context. We then now move to the cases of multiple environments.

7.2 Multiple environments (POMDP)

We will introduce some new notations of environment transitions. We note that after uniformization, the probability of transitioning from one environment to itself (the diagonal of table 2) does not matter any more since it will be taken into account in the dummy event. Furthermore, the sum of the other probability does not need to be less than 1 as we will divide them with the normalization constant.

The more important thing now is the ratio between each probabilities pair. We can rewrite table 2 as follow:

Table 11: Uniformized matrix transition of the environment of Q_i

	m_1	m_2	...	m_m
m_1	*	$\xi_{i,1,2}$...	$\xi_{i,1,m}$
m_2	$\xi_{i,2,1}$	*	...	$\xi_{i,2,m}$
\vdots				
m_m	$\xi_{i,m,1}$	$\xi_{i,m,2}$...	*

Furthermore, because we want to see if the speed of environment changing affects the optimal policy of each method, we multiply all transition matrices of each Q_i of the system with α . This does not change the ratio mentioned above and only takes effect when uniformization. In this process, it acts as a ratio between the change in the observable state and non observable (environment) state. If α is small, the rate of environment changing will only take a small part in the uniformized probability and vice versa.

From all the small transition matrices of each Q_i , we can get a big environment transition matrix $\mathcal{M}_{\mathcal{E}}$ where each row and column represents an environment state encoded using the method present in [From \$n\$ -dimensions to 2-dimensions](#) section. Since all the entries except the ones on the diagonal are multiplied by α , we can express $\mathcal{M}_{\mathcal{E}}$ as follow:

$$\mathcal{M}_{\mathcal{E}} = \begin{pmatrix} \rho_{1,1} & \alpha\rho_{1,2} & \dots & \alpha\rho_{1,|\mathcal{E}|} \\ \alpha\rho_{2,1} & \rho_{2,2} & \dots & \alpha\rho_{2,|\mathcal{E}|} \\ \vdots & & & \vdots \\ \alpha\rho_{|\mathcal{E}|,1} & \alpha\rho_{2,2} & \dots & \rho_{|\mathcal{E}|,|\mathcal{E}|} \end{pmatrix}$$

Where for example if $(m_1, m_1, \dots, m_1, m_1) \mapsto i$ and $(m_1, m_1, \dots, m_1, m_2) \mapsto j$, then $\rho_{i,j}$ will be equal to $\xi_{n,1,2}$ (the probability that the environment of queue n moves from m_1 to m_2).

If $\Pi_{\mathcal{E}}$ is the solution of equation (12) then $\mathcal{M}_{\mathcal{E}}\Pi_{\mathcal{E}} = \Pi_{\mathcal{E}}$ (because $P(E'|E)$ is actually ρ).

Lastly, with all the notions above, we can prove that the solution of equation (12) does not depend on α .

Lemma. *Given a set of all $\rho_{i,j}$ ($i \neq j$) defined as above.*

We call a pair of $(\alpha, (\rho_{1,1}, \dots, \rho_{|\mathcal{E}|,|\mathcal{E}|}))$ satisfied if \mathcal{M}_{α} is a right stochastic matrix where:

$$\mathcal{M}_{\alpha} = \begin{pmatrix} \rho_{1,1} & \alpha\rho_{1,2} & \dots & \alpha\rho_{1,|\mathcal{E}|} \\ \alpha\rho_{2,1} & \rho_{2,2} & \dots & \alpha\rho_{2,|\mathcal{E}|} \\ \vdots & & & \vdots \\ \alpha\rho_{|\mathcal{E}|,1} & \alpha\rho_{2,2} & \dots & \rho_{|\mathcal{E}|,|\mathcal{E}|} \end{pmatrix}$$

If $(\alpha, (\rho_{1,1}, \dots, \rho_{|\mathcal{E}|,|\mathcal{E}|}))$ and $(\alpha', (\rho'_{1,1}, \dots, \rho'_{|\mathcal{E}|,|\mathcal{E}|}))$ are satisfied, and there is T such that $\mathcal{M}_{\alpha}T = T$, then $\mathcal{M}_{\alpha'}T = T$

Proof. From $\mathcal{M}_{\alpha}T = T$, we have:

$$\begin{aligned} &\Leftrightarrow \alpha \sum_{j,j \neq i} \rho_{i,j} \tau_j + \rho_{i,i} \tau_i = \tau_i \\ &\Leftrightarrow \alpha \sum_{j,j \neq i} \rho_{i,j} \tau_j + (1 - \alpha \sum_{j,j \neq i} \rho_{i,j}) \tau_i = \tau_i \\ &\Leftrightarrow \alpha \sum_{j,j \neq i} \rho_{i,j} \tau_j - \alpha \tau_i \sum_{j,j \neq i} \rho_{i,j} = 0 \\ &\Leftrightarrow \sum_{j,j \neq i} \rho_{i,j} \tau_j = \tau_i \sum_{j,j \neq i} \rho_{i,j} \\ &\Leftrightarrow \alpha' \sum_{j,j \neq i} \rho_{i,j} \tau_j = \alpha' \tau_i \sum_{j,j \neq i} \rho_{i,j} \\ &\Leftrightarrow \alpha' \sum_{j,j \neq i} \rho_{i,j} \tau_j + (1 - \alpha' \sum_{j,j \neq i} \rho_{i,j}) \tau_i = \tau_i \\ &\Leftrightarrow \alpha' \sum_{j,j \neq i} \rho_{i,j} \tau_j + \rho'_{i,i} \tau_i = \tau_i \end{aligned}$$

Therefore, $\mathcal{M}_{\alpha'}T = T$. □

7.2.1 Scheduling

7.2.1.1 Small α

We start with a simple parameter as below:

Table 12: Scheduling 1

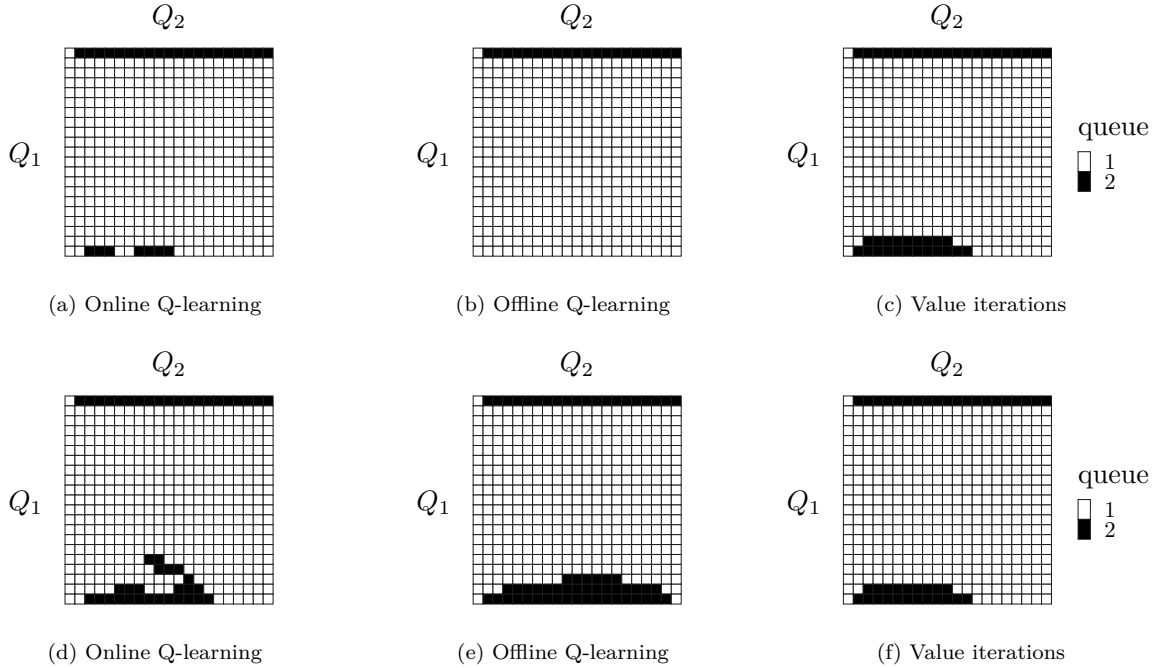
		m_1	m_2
C_1	Holding cost	5	1
	Arrival rate	0.13	0.13
	Departure rate	0.3	0.3
C_2	Holding cost	2	
	Arrival rate	0.13	
	Departure rate	0.3	

Table 13: Scheduling 2

		m_1	m_2
C_1	Holding cost	1	5
	Arrival rate	0.13	0.13
	Departure rate	0.3	0.3
C_2	Holding cost	2	
	Arrival rate	0.13	
	Departure rate	0.3	

Table 14: Matrix transition of the environment of Q_1 of table 12 and table 13

	m_1	m_2
m_1	*	1
m_2	1	*
α	10^{-3}	

Figure 16: Result of table 12 (upper) and 13 (lower) for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

First, for both cases, the solution of equation (12) is $\Pi = (0.5, 0.5)$ with $\mathcal{E} = \{(m_1, m_1), (m_2, m_1)\}$ (since queue 2 has only one environment). Furthermore, from equation (13), the value iteration method is actually a MDP problem as above with the following parameters.

As in [BVW85], the solution of table 15 is serving queue 1, which is also the result of the value iteration

Table 15: Empirical MDP of table 12 and table 13

Q_1	Holding cost c_1	$0.5 \times 1 + 0.5 \times 5 = 3$
	Arrival rate λ_1	$0.5 \times 0.13 + 0.5 \times 0.13 = 0.13$
	Departure rate μ_1	$0.5 \times 0.3 + 0.5 \times 0.3 = 0.3$
Q_2	Holding cost c_2	2
	Arrival rate λ_2	0.13
	Departure rate μ_2	0.3

method on both table 12 and table 13.

And finally, Q-learning also captures this information quite well and gives us a similar result instead of choosing which queue has less cost in the initial environment (queue 1 for table 12 and queue 2 for table 13) even though the speed of environment changing is slow.

We now move to a more complex case as follow:

Table 16: Scheduling 3

		m_1	m_2
C_1	Holding cost	3	4
	Arrival rate	0.13	0.12
	Departure rate	0.28	0.31
C_2	Holding cost	2	6
	Arrival rate	0.13	0.14
	Departure rate	0.3	0.32

Table 17: Matrix transition of the environment of Q_1 of table 16

	m_1	m_2
m_1	*	1
m_2	2	*
α	10^{-3}	

Table 18: Matrix transition of the environment of Q_2 of table 16

	m_1	m_2
m_1	*	4
m_2	3	*
α	10^{-3}	

Using the same logic as above, we can transform it into a MDP problem as follow:

Table 19: Empirical MDP of table 16

Q_1	Holding cost c_1	3.3333
	Arrival rate λ_1	0.1266
	Departure rate μ_1	0.2899
Q_2	Holding cost c_2	4.2857
	Arrival rate λ_2	0.1357
	Departure rate μ_2	0.3114

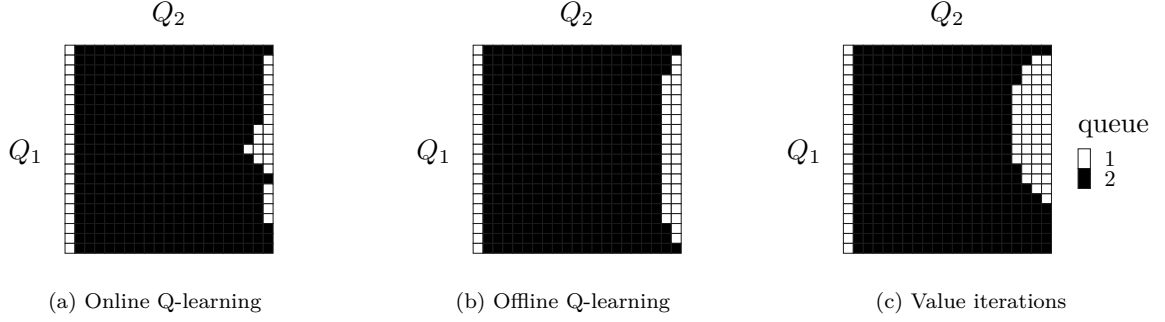


Figure 17: Result of table 16 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

The result of value iteration on the right side is also the same as the value iteration of table 19. Both actually gives the same numerical estimation of V_* . Q-learning is also good at learning this type of system and we can see that it deduces a similar policy to the value iteration's one.

7.2.1.2 Big α

We take the same parameters as in table 13 with a different α .

Table 20: Matrix transition of the environment of Q_1 of table 13

	m_1	m_2
m_1	*	1
m_2	1	*
α	10^2	

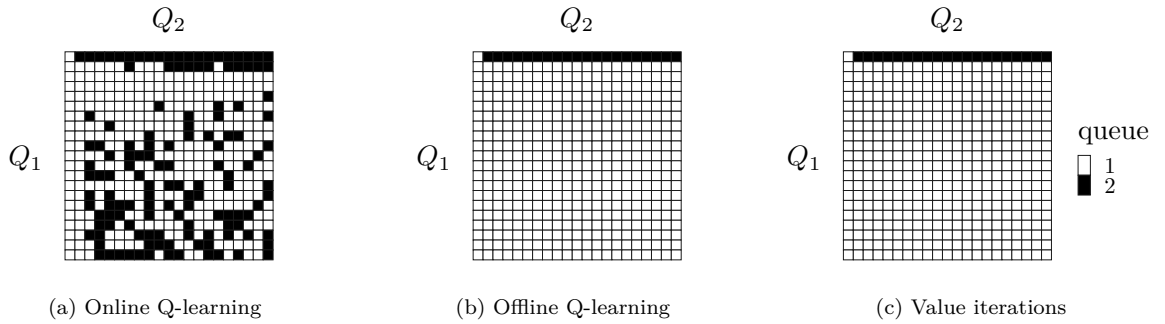


Figure 18: Result of table 20 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

As explained above, the policy perceived by the value iteration is quite similar as above. The majority of online Q-learning cells are choosing action 1, especially the ones near the origin. However, Q-learning still contains a lot of noise, a sign that suggests we have not iterated enough. Although the number of iterations (3×10^{10}) is sufficiently large, bigger α requires a lot more iterations because in that case, the system does not evolve much in trajectory but in environment instead (for example we could think that the system will stuck at one cell near the origin and the environment state constantly changes after each iteration). The figure below shown that phenomena:

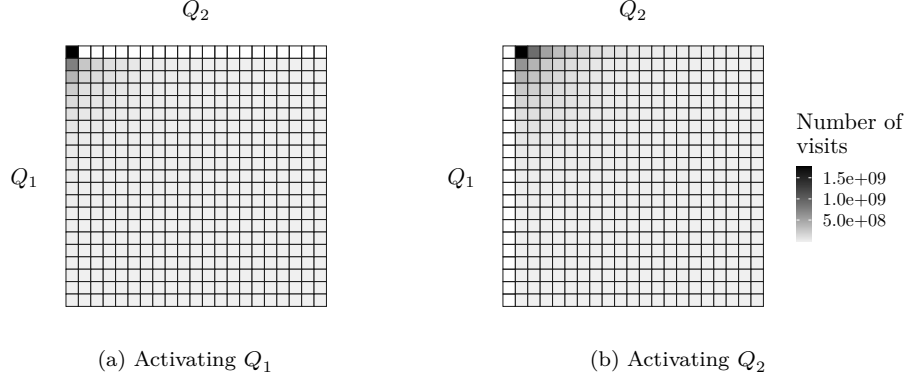


Figure 19: Number of visits of online Q-learning on table 20 for 20x20 case. The highest number of visits is $\approx 4 \times 10^9$ (at $(0, 0)$) while the lowest one is $\approx 5 \times 10^5$ (at $(20, 20)$).

For offline Q-learning, we see no difference with the value iteration method. As in the explanation for the result of figure 11, in the same amount of time, most of state-action pairs in offline version receives more visits than the online one (except the cells around the origin) and therefore gives a less noisy result.

We also try the complex example of table 16 with $\alpha = 10^2$. Below is the policy perceived by Q-learning and value iteration.

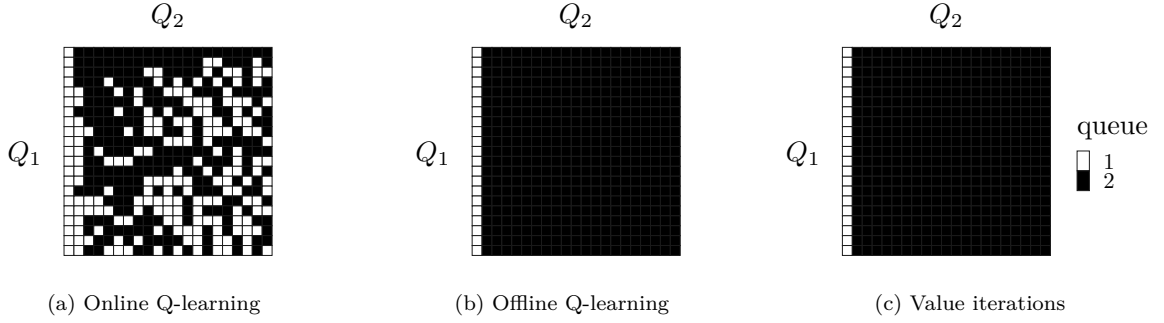


Figure 20: Result of table 16 with $\alpha = 10^2$ for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

Same as above, offline Q-learning and value iteration method give the same policy while the majority of Q-learning choose the optimal action of the empirical MDP. But it still contains some noisy cells that have not received enough visits.

7.2.2 Load-balancing

7.2.2.1 Small α

We will start with a simple case as follow:

Table 21: Load-balancing 1

		m_1	m_2
Q_1	Holding cost	3	1
	Departure rate	0.3	0.3
Q_n	Holding cost	2	
	Departure rate	0.3	
Global arrival rate		0.5	

Table 22: Matrix transition of the environment of Q_1 of table 21

	m_1	m_2
m_1	*	1
m_2	1	*
α	10^{-3}	

Table 23: Empirical MDP of table 21

Q_1	Holding cost	2
	Departure rate	0.3
Q_n	Holding cost	2
	Departure rate	0.3
	Global arrival rate	0.5

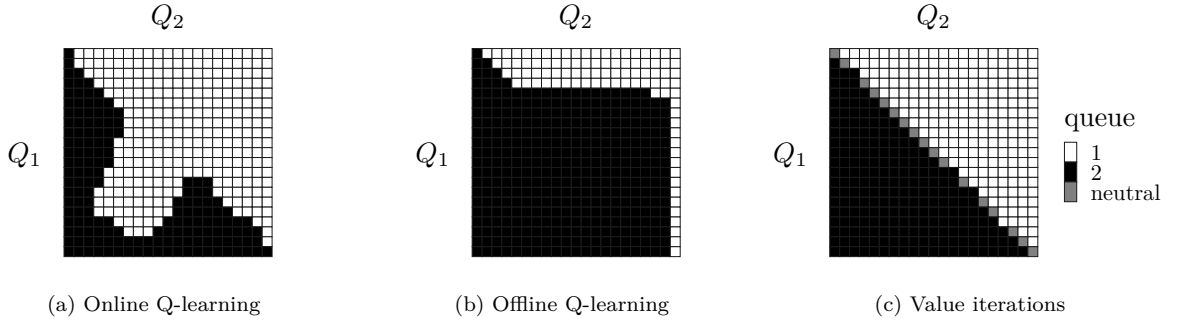


Figure 21: Result of table 21 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

The value iteration method shows the same result as expected with a neutral diagonal line. For online Q-learning, Although there is a large hole in the center of the grid, we could still see a diagonal line that split the grid into two regions as in the optimal policy of the empirical MDP. However, the offline method gives a very different result compared to the other methods.

Since the learning capacity of Q-learning is not really good for complex load balancing systems, we will just test it with the same case but bigger α .

7.2.2.2 Big α

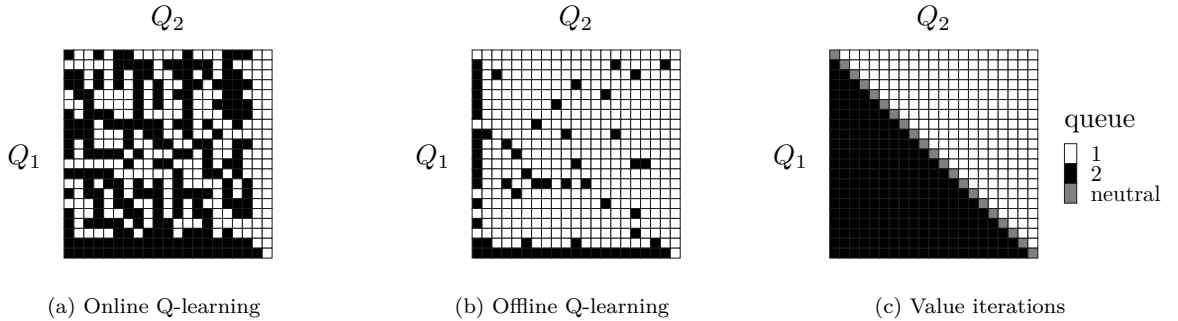


Figure 22: Result of table 21 for 20x20 case. $\gamma = 0.999$, $\epsilon = 0.999$.

Unfortunately, the online Q-learning result here is completely random and we will need to run a much longer simulation for it to be converged. With offline version, it gives a very different, nearly opposite

policy with the policy for small α . We can see that although the offline method works well with the scheduling system, it is broken on the load-balancing system.

8 Conclusion and extension

First, the table below summarizes the advantages and weaknesses of each method in our report:

Table 24: Comparison between methods for MDP / POMDP problem.

Aspect	Online Q-learning	Offline Q-learning	Value iterations
Can be used on POMDP	Yes	Yes	No
Prerequisites	A system simulator or nothing (by interacting directly with the real system)	A reliable system simulator	Real or calculated (empirical) parameters of the system such that it is a MDP problem
Checking for convergence (checking if the value does not change much after some iterations)	Since only one cell is updated at a time, and each cell is not visited uniformly, such checking requires a lot of efforts and will slow down the simulation. An alternative solution is re-running the simulation with different learning steps and comparing the result, but this also requires a lot of time and manual effort.	By calculating $\ \cdot\ _2$ between the previous Q-value and the current Q-value.	By calculating $\ \cdot\ _2$ between the previous V-value and the current V-value.
Speed of convergence	Since it follows the trajectory, one training session might take a very long time. In addition, because checking for convergence is not simple, we might have to run the training session multiple times as well as carefully tuning the hyper-parameters so that takes a lot of time.	Faster than online Q-learning but still taking long time before converging	The fastest method.
Quality of the result	The result might not converge to the best policy and contain a lot of noise. Furthermore, its result is extremely sensitive to system parameters, especially the ones having impacts on the evolution of the whole system.	The result might not converge to the best policy but it will less likely to contain noise and is more robust to a change in the settings of system.	The result converged to the optimal solution when we are working on a MDP problem.

Aspect	Online Q-learning	Offline Q-learning	Value iterations
Behavior on POMDP	Online Q-learning seems to perceive the problem as an empirical MDP problem as we could observe that its results resemble the value iteration result. But we need more analytical inspections to verify this hypothesis.	Same as online Q-learning	This method can not work on a POMDP problem, but said problem can be converted to an empirical MDP problem by the method describe in section Empirical transition probability .

Through the report, we have seen different behaviors of each method in solving MDP and POMDP. In conclusion, Q-learning can learn the underlying structure of simple systems (as seen in figure [8](#), [9](#), [10](#), [16](#), [17](#)). But when the system becomes more complex, it might take much longer to capture this information (the case of figure [18](#), [21](#)) or completely fail to do so (figure [11](#), [13](#), [14](#), [15](#), [22](#)). Furthermore, modeling and training takes a particularly long time so more work is needed to have a complete understanding about what Q-learning learns in a system with partial information.

Finally, we will end the report with an unanswered question: **“How can we know that the policy we obtained is the optimal policy in case of POMDP ?”**. This question is very important because if we have the optimal policy beforehand, not necessarily through a model free method, we can later compare it with the result of another method to get an idea if that method works with our problem or not (like what we did with the MDP problem: compare the Q-learning result and the value iteration one). However, because of the limited context of the internship, we are not able to answer that question and left it as a possibility for future work.

References

- [BVW85] C. Buyukkoc, P. Varaiya, and J. Walrand. “The $c\mu$ rule revisited”. In: *Advances in Applied Probability* 17.1 (1985), pp. 237–238. DOI: [10.2307/1427064](https://doi.org/10.2307/1427064).
- [WD92] Christopher J. C. H. Watkins and Peter Dayan. “Q-learning”. In: *Machine Learning* 8.3 (May 1992), pp. 279–292. ISSN: 1573-0565. DOI: [10.1007/BF00992698](https://doi.org/10.1007/BF00992698). URL: <https://doi.org/10.1007/BF00992698>.