

Privacy Preserving Q-learning in the Analog Model for Secure Multiparty Computation

Hirofumi Miyajima^{†1}, Noritaka Shigei^{†2}, Hiromi Miyajima^{†3} and Norio Shiratori^{†4}

Abstract—Many studies have been done with the security of cloud computing. Though data encryption is a typical approach, high computing complexity for encryption and decryption of data is needed. Therefore, safe system for distributed processing with secure data attracts attention, and a lot of studies have been done with them. SMC (Secure Multiparty Computation) is one of these methods. So far, most of works for ML (Machine Learning) with SMC are ones with supervised and unsupervised learning such as BP (Back Propagation) and K-means methods. Then, in the case where learning data does not exist explicitly like reinforcement learning (RL), how should it be done? We already proposed some algorithms of Q-learning and PS (Profit Sharing) learning for SMC in previous papers. However, they were methods using digital models. On the other hand, solutions for analog models are desired in the real world. In this paper, we propose SMC algorithms for Q-learning in the analog model and show their effectiveness. The idea is that in the digital model, only one behavior is selected at each time, whereas in the analog model it is decided as a combination of a plural of weighted actions.

Index Terms—cloud computing, secure multiparty computation, Q-learning, analog model.

I. INTRODUCTION

WITH increasing interest in Artificial Intelligence (AI), many studies have been made with Machine Learning (ML). With ML, the supervised, the unsupervised and RL are well known. Recently, the importance of RL in AI is increasing with the advancement of learning. Due to respond to the increase in the amount of data or complex problems for ML, the use of cloud computing systems is spreading. The development of cloud computing allows the use such as big data analysis to analyze enormous information accumulated by the client, and to create market value of data [1]–[6]. On the other hand, the client of cloud computing system cannot escape from anxiety about the possibility of information being abused or leaked. In order to solve the problem, data processing methods can be considered such as cryptographic one [1], [2]. However, data encryption system requires both encryption and decryption for requests of client or user, so its applications are limited. Therefore, safe systems for distributed processing with secure data attract attention, and a lot of studies with them have been done. It is known

that SMC's idea of distributing learning data among multiple servers is one method to realize this [3], [4]. As for SMC, many methods of learning by sharing learning data into partial subsets have been proposed [3]–[6]. Then, in the case where learning data does not exist explicitly like RL, how should it be done? We have proposed an applicable model also in this case. The idea is to divide not only learning data but also learning parameters, find partial solutions at each server, combine them and make it the solution of the system [8]. Based on this idea, we already proposed some algorithms of Q-learning and PS learning for SMC in previous papers [9], [10]. However, they were methods using the digital model. On the other hand, solutions to analog model are desired in the real world [11], [12]. The latter gives a solution closer to the optimal solution than the former.

In this paper, we propose SMC algorithms for Q-learning in the analog model and show their effectiveness. The idea is that in the digital model, only one action is selected at each time, whereas in the analog model it is decided as a combination of weighted actions. In Section 2, we explain cloud computing system, related works on SMC and how to share the data used in this paper. Further, a Q-learning method in the analog model is introduced. In Section 3, we propose two Q-learning methods in the analog model for SMC. In section 4, numerical simulations for a maze problem are performed to show the performance of proposed methods.

II. PRELIMINARY

A. Cloud system and related works with SMC

The system used in this paper is composed of a client and m servers. Each data is divided into m pieces of numbers and is sent to each server (Fig.1 for $m = 2$). Each server performs its computation and sends the computation result to the client. The client can get the result using them. If the result is not obtained by one processing, then the multiple processing is repeated. As for the cloud system, there are many methods of secure preserving, but it seems that SMC method using distributed processing is suitable for the system. In particular, three types of conventional methods for partitioning data to be securely shared are well known [3], [4]. They are known as horizontal, vertical and arbitrary partitioning methods. In the following, the horizontal method is only explained by using a data example of students' marks shown in Table I. See Miyajima [8] about the detailed explanation. In Table I, a and b are original data (marks) and ID is the identifier of student. The number of servers is two. The assumed task is to calculate the average of data. The horizontal partitioning method assigns the horizontally partitioned data to servers as follows :

Server 1 : data for ID = 1, 2

Affiliation: Faculty of Informatics, Okayama University of Science, 1-1 Ridaicho, Kitaku, Okayama, 700-0005, Japan

corresponding author to provide email: miya@mis.ous.ac.jp

^{†1} email: miya@mis.ous.ac.jp

Affiliation: Kagoshima University, Kagoshima, Japan

^{†2} email: shigei@eee.kagoshima-u.ac.jp

^{†3} email: miya@eee.kagoshima-u.ac.jp

Affiliation: Research and Development Initiative, Chuo University, Tokyo, Japan

^{†4} email: norio@shiratori.riec.tohoku.ac.jp

This work was supported in part by the JSPS KAKENHI GRANT NUMBER JP17K00170.

TABLE I
 CONCEPT OF HORIZONTALLY AND VERTICALLY PARTITIONED METHODS
 COMPOSED OF ONE CLIENT AND TWO SERVERS.

	ID	Subject A	Subject B	
		a	b	
Server 1	1	22	32	Horizontally partitioned method
	2	24	37	
Server 2	3	40	40	
	4	13	45	
Average		24.75	38.5	Vertically partitioned method

Server 2 : data for ID = 3, 4

In the method, Server 1 computes two averages for subjects A and B as $(22 + 24)/2$ and $(32 + 37)/2$, respectively. Likewise, Server 2 computes two averages for subjects A and B as $(40 + 13)/2$ and $(40 + 45)/2$, respectively. Servers 1 and 2 send calculated averages to the client and the client obtains the averages of subjects A and B as 24.75 and 38.5, respectively. Since each server cannot know half of the dataset, the method preserves privacy (See Table I). Remark that each server also cannot know the result when we consider average values as unknown parameters. The vertical partitioning method is one of processing data for each subject (See Table I). The third method, the arbitrary partitioning method, splits horizontally and vertically the dataset into multiple parts, and the method assigns the split parts to the servers. For any of the above mentioned methods, if the number of servers is fewer, that is, the size of a partitioned data is larger, any server may more easily guess the feature of all the data from its own subset of data. Therefore, the methods need a large number of servers in order to keep privacy and security. On the other hand, the method explained in the next section divides each item of data and seems to keep them even in the case of a small number of servers.

B. Secure divided data for SMC and their application to machine learning

Let us explain secure divided data for the proposed method using Table II [8]. Let a and b be two marks and $m = 2$ (See Fig.1). Assume that the addition form is used for dividing each item. For example, two marks a and b are divided into two real numbers as $a = a_1 + a_2$ and $b = b_1 + b_2$ as follows:
 $a = a_1 + a_2 : a_1 = 1(r_1/10), a_2 = a(1 - r_1/10)$
 $b = b_1 + b_2 : b_1 = b(r_1/10)$ and $b_2 = b(1 - r_1/10)$
 where r_1 is a real random number for $-9 \leq r_1 \leq 9$. If $r_1 = -1$, then $a_1 = 0.2$ and $a_2 = -2.2$ are obtained. Remark that Server 1 and Server 2 have all the data in column-wise of a_1 and b_1 and a_2 and b_2 for each ID as shown in Table II, respectively.

Let us explain how to compute the average for subject A using data a . Server 1 and Server 2 compute the average of a_1 and a_2 , respectively. In this case, each average in column-wise for a_1 and a_2 is 1.8 and -4.8 , respectively. As a result, the average is obtained as -3 from $1.8 - 4.8$.

Remark that each data for server is randomized and the method does not need to use complicated computation as the encryption system.

Let us explain about the application of secure divided data to ML. In the conventional methods, the set of learning data

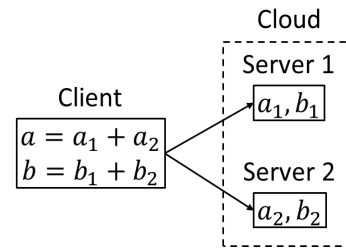


Fig. 1. An example of secure shared data for $m = 2$.

TABLE II
 DATA FOR SERVER 1 AND SERVER 2.

ID	subject A a	subject B b	Additional form				
			r_1	a		b	
				a_1	a_2	b_1	b_2
1	-2	-6	-1	0.2	-2.2	0.6	-6.6
2	-8	2	-6	4.8	-12.8	-1.2	3.2
3	1	-9	5	0.5	0.5	-4.5	-4.5
Server 1				1.8		-1.7	
Server 2					-4.8		-2.6
average	-3	-4.3					

for ML is shared into some subsets. On the other hand, the proposed one divides each item of the learning data into plural pieces and processes them. From the point of view, SMC algorithms for supervised learning such as BP method and unsupervised learning like k-means method were proposed [8]. Then, how is the algorithm of RL for SMC? In this case, as there do not exist learning data explicitly, the solution is obtained by repeating trial and error. Since there is no data for learning, it seems that the conventional method using subset of learning data is difficult to use. On the other hand, several methods on privacy preserving for RL have been proposed, but these are almost all methods using encryption and homomorphic mapping [13], [14]. The proposed method attempts to realize SMC by simple secret computation processing which does not use such complicated cryptographic processing and homomorphic mapping. That is, the aim to reduce the computational complexity of client while keeping the secret of data (Q-value information in this case). In previous papers, we have already proposed Q-learning and PS methods for SMC in the digital model [9], [10]. In the next section, learning methods using secure divided data in the analog model are proposed.

C. Q-learning methods in the digital and analog models

First, let us explain about the Q-learning algorithm in the digital model. The Q-learning is one of RL techniques for environment-identity type [11]. It can be used to find an optimal action-selection policy for a given Markov Decision Process (MDP). In solving problems using Q-learning, it is determined how the agent selects an action at any state. It is performed by learning an action-value (Q-value) function that gives the expected utility of taking the action for the current state [11], [12]. The Q-value function is defined as a function $Q : S \times A \rightarrow R$, where S , A and R are sets of states, actions and real numbers, respectively. First, let all Q-values be 0. Then each action for a state is selected randomly. Boltzmann or ϵ -greedy method is used as the method selected data randomly as shown later. If a desired solution for the problem is not obtained, learning is iterated.

In the Q-learning, the action is selected based on the Q-value function. If a solution is obtained, Q-values are updated based on the updated formula. By iterating these process, it is known that Q-value function is updated and converges [11]. In the first part of learning, the action for the state is selected randomly and the action becomes decidable as learning steps proceed.

In learning step, Q-value function is updated as

$$Q(s, a) \leftarrow Q(s, a) + \alpha \Delta \quad (1)$$

$$\Delta = r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a) \quad (2)$$

where r , α and γ are the reward, learning constant and discount rate, respectively. The state s' is the next state selected for the state s and the action a . The term $\max_{a' \in A} Q(s', a')$ means the Q-value $Q(s', a'_0)$ for an action a'_0 taking the maximum number of $Q(s', a')$. See the Ref. [9] about the example of Q-learning in the digital model.

Let us think about the analog model of Q-learning. Let us explain using an example of maze problem to make the story easy to understand. The problem is how the agent can arrive in the shortest path at the goal from the start point (See Fig.2). In the digital model, one of actions at each position is selected. Therefore, the path obtained becomes a zigzag path as behavior. However, in the real problem, it is desirable to get smooth path. An analog model is proposed as a model to realize all directions for behavior (action). Several methods have been proposed for Q-learning for the analog model. Here, we introduce a learning method that determines the action based on the distance between the current position and state (position) of the agent. That is, while hard selection for the action in the digital model is performed, the introduced method achieves a soft matching selection for the action. Let us explain how to select the action of Q-learning in the analog model using Fig.2. Assign n states (the center position c) in the space where the agent moves.

Let d be the current position of the agent. First, the distance $D_j(d)$ between d and the center position c_j of each state for $1 \leq j \leq n$ is computed as follows:

$$D_j(d) = \exp\left(-\frac{\|d - c_j\|^2}{2b_j^2}\right) \quad (3)$$

where $\exp(\cdot)$ means the Gauss function with the width b_j . Remark that $D_j(d)$ is large if the distance between d and c_j is near.

Next, the action a_j^* for each state s_j is selected as the action with the maximum number of Q-value (See Fig.2) and the degree μ_j of coincide between d and s_j is computed as follows:

$$\mu_j(d) = \frac{D_j(d)}{\sum_{i=1}^n D_i(d)} \text{ for } 1 \leq j \leq n \quad (4)$$

and

$$\sum_{j=1}^n \mu_j(d) = 1 \quad (5)$$

Further, the action a^* for the agent is determined as the composition of vectors as follow:

$$a^* = \sum_{j=1}^n \mu_j a_j^* \quad (6)$$

$s \setminus d$	Distance D	Degree μ	a					$\max_i Q_j$
			a_1	\dots	a_i	\dots	a_M	
$s_1(c_1)$	$D_1(d)$	μ_1	Q_{11}	\dots	Q_{1i}	\dots	Q_{1M}	$\rightarrow a_1^*$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$s_j(c_j)$	$D_j(d)$	μ_j	Q_{j1}	\dots	Q_{ji}	\dots	Q_{jM}	$\rightarrow a_j^*$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$s_n(c_n)$	$D_n(d)$	μ_n	Q_{n1}	\dots	Q_{ni}	\dots	Q_{nM}	$\rightarrow a_n^*$

$\Rightarrow a = \sum_{j=1}^n \mu_j a_j^*$

Fig. 2. The method to determine the action in the analog model.

As a result, the agent at the position d moves in the direction of a^* and arrives at the new position d' . In this case, the moving distance is selected randomly.

The conventional algorithm for Q-learning is shown as follows :

[Q-learning algorithm for the analog model]

d : the current position.

d_0 : the initial position.

$Q(s, a)$: Q-value for the state s and the action a .

S : The set of states.

A : The set of actions.

t_{max} : The maximum number of learning time.

T_{max} and T_{min} : Constants for Boltzmann selection.

c_j : The center position of the state s_j for $1 \leq j \leq n$.

$D_j(d)$: The distance between the state s_j and the place d for Eq.3.

ε : The probability for ε -greedy selection method.

p : The real number selected randomly from the interval $0 \leq p \leq 1$.

Step 1

Let R^+ and R^- , α and γ be plus and minus reward, learning constant and discount rate. Let S and A be defined. Let $Q(s, a) = 0$ for $s \in S$ and $a \in A$. Let $t = 0$. (Let ε be defined for ε -greedy selection method.)

Step 2

Let $d \leftarrow d_0$ be the start position.

Step 3

The distance $D_j(d)$ and the degree μ_j for $1 \leq j \leq n$ are computed from Eqs.3 and 4.

Step 4 (In the case of Boltzmann selection method)

The action a_j^* at the state s_j is selected based on the value $B(s_j, a)$ as follows :

$$B(s_j, a) = \frac{\exp(Q(s_j, a)/T)}{\sum_{b \in A} \exp(Q(s_j, b)/T)} \quad (7)$$

$$T = T_{max} \times \left(\frac{T_{min}}{T_{max}}\right)^{\frac{t}{T_{max}}} \quad (8)$$

Let a_j^* be the selected action based on Eqs.7 and 8.

That is, the action a_j^* from the set A is selected with the probability $B(s_j, a)$.

Step 5

The action a^* at the place d is computed using Eq.6.

Let d' be the next position determined from the action a^* .

Step 6

Each value of $Q(s_j, a_j^*)$ for the state s_j and the action a_j^* is updated as follows :

$$Q(s_j, \mathbf{a}_j^*) \leftarrow Q(s_j, \mathbf{a}_j^*) + \mu_j \alpha (r + \Delta) \quad (9)$$

$$\Delta = \sum_{j=1}^n \{ \gamma \mu_j' \max_{a_j' \in A} Q(s_j', a_j') - \mu_j Q(s_j, a_j^*) \} \quad (10)$$

, where μ_j' , s_j' and a_j' are parameters for \mathbf{d}' and if \mathbf{d}' is in the goal area, then $r = R^+$ and go to Step 7 else if \mathbf{d}' is not in the movable place(state), then $r = R^-$ and go to Step 3 else $r = 0$ and go to Step 3.

Step 7

If $t = t_{max}$ then the algorithm terminates else go to Step 3 with $t \leftarrow t + 1$.

If the ε -greedy method is used instead of Boltzmann selection, Step 4 is replaced with the following Step 4' :

Step 4'(In the case of the ε -greedy selection method)

Let p be the real number randomly selected. The action a_j^* for the state s_j is selected based on ε -greedy selection as follows :

$$\mathbf{a}_j^* = \begin{cases} \mathbf{a} \text{ s.t. } \max_{\mathbf{a} \in A} Q(s_j, \mathbf{a}) & \text{for } p \geq \varepsilon \quad (a) \\ \text{randomly selected} & \text{for otherwise} \quad (b) \end{cases} \quad (11)$$

That is, if p is greater than or equal to ε , then the action \mathbf{a}_j^* satisfying the condition of Eq.11(a) is selected otherwise the action \mathbf{a}_j^* is randomly selected.

III. Q-LEARNING FOR SECURE MULTIPARTY COMPUTATION IN THE ANALOG MODEL

In Q-learning for SMC on cloud system, Q-values are divided to each server in addition form . Each server updates divided Q-values and sends the result to the client. The client can get new Q-values by adding the results of m servers. The process is iterated until the evaluating value for the problem satisfies the final condition. The problem is how Q-values on the client are updated using Q-values divided on each server. The divided representation of Q-value is given as follows :

$$Q(s, a) = \sum_{k=1}^m Q_k(s, a) \text{ for } s \in S \text{ and } a \in A \quad (12)$$

In the following, two learning methods are proposed.

The proposed methods can be easily applied to other Q-learning algorithms in the analog model.

A. Q-learning for SMC

The first algorithm using Boltzmann (ε -greedy) selection method is shown in Table III. The initial values of client and servers are set in Initializing Step. In Step 1, the initial position is set. In Step 2, the distance $D_j(\mathbf{d})$ and the degree μ_i for the current position \mathbf{d} and each center position c_j of state s_j for $1 \leq j \leq n$ are computed. In Step 3, the action \mathbf{a}_j^* for the state s_j is selected based on Boltzmann selection method. Further, the next action \mathbf{a}^* is determined as the vector sum of \mathbf{a}_j^* for $1 \leq j \leq n$ and the new position \mathbf{d}' is obtained from the action \mathbf{a}^* . Furthermore, the degree μ_j' of coincide for the position \mathbf{d}' is computed from \mathbf{d}' and s_j ($1 \leq j \leq n$). In Step 4, the updating rate Δ_b^k for $b \in A$ and $1 \leq k \leq m$ is computed

in each server. In Step 5, the updating rate Δ^* is computed using Δ_b^k for $b \in A$ and $1 \leq k \leq m$. Further, the updating rate ξ_k for each server is determined using $\beta_k(s_j, \mathbf{a}_j^*)$ and Δ^* . It means to divide Δ^* into m pieces of numbers. In Step 6, the Q-value $Q_k(s_j, \mathbf{a}_j^*)$ for each server is updated. In Step 7 and Step 8, it is checked if the final condition satisfied or not. If the final condition is not satisfied, then the next episode is iterated. The algorithm is called M1 method.

B. Q-learning with dummy updating

In M1 method, designated Q-values are updated based on μ_j for $1 \leq j \leq n$. Therefore, there is the possibility that the server knows information on which Q-value is important. In order not to inform the server of the update information, an improved method with dummy updating is proposed in which all the Q-values are updated.

The fundamental idea of Q-learning with dummy updating is that all Q-values are updated at each step. Therefore, it seems that each server cannot know which Q-value is important or not. Let us explain the improved M1 method. The number $p_k(s, \mathbf{a})$ is randomly selected such that $|p_k(s, \mathbf{a})| \leq 1$ and $\eta_k(s, \mathbf{a})$ is calculated as follows :

$$\eta_k(s, \mathbf{a}) = \begin{cases} \frac{p_k(s, \mathbf{a})}{\sum_{l=1}^m p_l(s, \mathbf{a})} & \text{for } s = s_j \text{ and } \mathbf{a} = \mathbf{a}_j^* \quad (A) \\ \frac{p_k(s, \mathbf{a})}{\sum_{l=1}^m p_l(s, \mathbf{a})} - \frac{1}{m} & \text{for } 1 \leq j \leq n \quad (B) \\ \frac{p_k(s, \mathbf{a})}{\sum_{l=1}^m p_l(s, \mathbf{a})} - \frac{1}{m} & \text{for otherwise} \quad (B) \end{cases} \quad (13)$$

where $s = s_j$ and $a = a_j^*$ for $1 \leq j \leq n$ are selected states and actions.

Remark that each case of A and B for Eq.13 holds $\sum_{k=1}^m \eta_k(s, \mathbf{a}) = 1$ and 0 for $s \in S$ and $\mathbf{a} \in A$, respectively. That is, all Q-values for states and actions are randomly updated using $p_k(s, \mathbf{a})$ and each server cannot know which Q-value is import or not at each step.

In the improved method, Steps 8 and 9 of Table III are changed as follows :

Step 8 (Client)

Calculate $\eta_k(s, \mathbf{a})$ of Eq.13 for $s \in S$ and $\mathbf{a} \in A$. Let $\xi_k = \alpha \eta_k(s, \mathbf{a}) \Delta^*$ for $1 \leq k \leq m$. Send ξ_k for $1 \leq k \leq m$ to each server.

Step 9 (k-th Server)

The Q-value $Q_k(s_j, a_i)$ is updated as follows :

$$Q_k(s_j, a_i) \leftarrow Q_k(s_j, a_i) + \mu_j \xi_k \quad (14)$$

The algorithm is called M2 method.

IV. NUMERICAL SIMULATIONS FOR THE PROPOSED ALGORITHMS

In numerical simulations, the problem is to find the shortest path for the agent from the start position to goal area by Q-learning methods (See Fig.3). In Fig.3, the agent can go to any position except for black and outer areas. In order to find the shortest path, the agent iterates trials to move from the start to goal area based on each algorithm. The simulation conditions are as follows:

1) Let the start position $\mathbf{d}_0 = (1, 1)$, black(wall) area $B = \{(x_1, x_2) \in R^2 | 4 \leq x_1, x_2 \leq 6\}$ and the goal area $\{(x_1, x_2) \in R^2 | 9 \leq x_1, x_2 \leq 10\}$, where R is the set of all real numbers.

TABLE III
 M1 METHOD OF Q-LEARNING FOR SMC.

	Client	k -th Server ($1 \leq k \leq m$)
Initialization	The numbers R^+ , R^- , r , α and γ are given. Let $t = 0$. Let $Q_k(s, a) = 0$ for $s \in S$ and $a \in A$.	
Step 1	Let $d \leftarrow d_0$.	
Step 2	Calculate $D_j(d)$ and μ_j for d and $c_j (1 \leq j \leq n)$.	
Step 3		Send all Q-values $Q_k(s, a)$ for $s \in S$ and $a \in A$ to the client.
Step 4	Calculate $Q(s, a) = \sum_{k=1}^m Q_k(s, a)$ for $s \in S$ and $a \in A$. Select the action s_j^* for the state s_j based on Boltzmann selection of Eq.(7). Let $a^* = \sum_{j=1}^n \mu_j a_j^*$ be the vector sum of $a_j^* (1 \leq j \leq n)$. Let d' be the next position determined by a^* . The degree μ_j' of coincide between d' and s_j for $1 \leq j \leq n$ is computed.	
Step 5	If the position d' is possible(movable), then send the degree μ_j' and μ_j for $1 \leq j \leq n$ to each server else go to Step 4.	
Step 6		Calculate $\Delta_b^k = r\mu_j' Q_k(s_j', b) - \mu_j Q_k(s_j, a_j^*)$ for $1 \leq j \leq n$ and send them to client.
Step 7	Calculate $\Delta_b = \sum_{k=1}^m \Delta_b^k$ and $\Delta^* = r + \max_{b \in A} \Delta_b$ and send $r + \Delta^*$ to all servers, where $r = R^+$, R^- and 0 if d' is in goal, not in movable position and otherwise, respectively.	
Step 8	Select m pieces of random numbers $\beta_k(s, a)$ s.t. $\sum_{k=1}^m \beta_k(s, a) = 1$. Let $\xi_k = \alpha \beta_k(s_j, a_j) \Delta^*$ for $1 \leq k \leq m$. Send ξ_k for $1 \leq k \leq m$ to each server.	
Step 9		The Q-value $Q_k(s_j, a_j^*)$ is updated as follows : $Q_k(s_j, a_j^*) \leftarrow Q_k(s_j, a_j^*) + \mu_j \xi_k$
Step 10	If d' is in goal state, then go to Step 11 else go to Step 2 with $d \leftarrow d'$.	
Step 11	If $t = t_{max}$ then the algorithm terminates else go to Step 1 with $t \leftarrow t + 1$.	

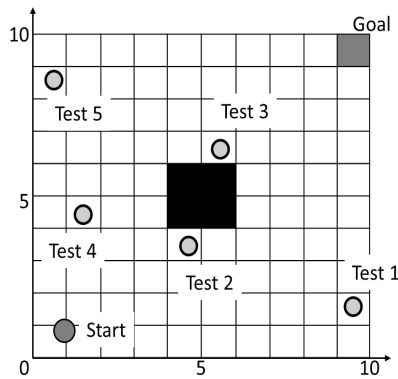


Fig. 3. The figure for the maze problem.

- 2) Let $A = \{0, \frac{\pi}{2}, \pi, \frac{3\pi}{2}\}$ and the set of central positions $C(S) = \{(x_1, x_2) \in Z^2 | 1 \leq x_1, x_2 \leq 9\} - B$, where each of the set A means four directions up, down, right and left on Fig.3, and Z is the set of all integers. Let $n = |C(S)|$. The states are arranged in a lattice pattern at regular intervals.
- 3) Let $T_{max} = 5.0$ and $T_{min} = 0.03$ for Boltzman selection and $\varepsilon = 0.1$ for ε -greedy selection.
- 4) The moving distance h for the agent at each position d is selected randomly for $0.01 \leq h \leq 0.5$ and the action is determined by each algorithm.
- 5) If the agent selects to move to wall or outer area, the agent ignores the selection and reselects a new action, where $r = -1$ is given as the negative reward. It is not counted in the number of trials.
- 6) If the agent arrives at the goal area in the maximum number of learning time, then the agent starts the new trial, where the reward $r = 1$ is given as the positive reward.
- 7) Let $t_{max} = 10000$, $r = 10$, $\alpha = 0.5$, $\gamma = 0.92$, and $m = 3$ and 10. Twenty trials for learning and test are performed for each algorithm.
- 8) In the test simulation, experiments are carried out with

five positions as the starting points as shown in Fig.3. The result is evaluated as the rate of the number of success trials and the average of moving distance from each starting point to goal area.

Figs.4 and 5 show the efficiency graphs. They represent the moving distance to the learning time. In Figs.4 and 5, the conventional (ε -greedy or Boltzmann), M1 for $m = 3$ and 10, and, M2 for $m = 3$ and 10 mean the conventional algorithm with ε -greedy or Boltzmann as selection method, proposed algorithms for M1 with $m = 3$ and 10, and for M2 with $m = 3$ and 10, respectively. All the results are almost the same as the conventional cases. Table IV shows the results of the test simulation, where Tests 1 to 5 mean the cases with different starting points as shown in Fig.3. In Table IV, No. Suc. and M.D. mean the number of success trials for twenty trials and the average of moving distance for success trials. Further, the result on each server means one in the cases where the same trials are performed using only Q-values of each server. For example, 0.95 of No.success and 9.83 of M.D. for Test 1 of client in the conventional method mean 19/20 and 9.83 steps as the average for 20 trials to the goal from the start point of Test 1. All the results for client are almost the same as the conventional. As for servers, trials only using server's information are almost unsuccessful and a lot of times are needed even in the successful cases.

Finally, let us explain about the result of selection methods. The ε -greedy method is in learning time faster than Boltzmann method. But, Boltzmann method is superior in accuracy of the test simulation to ε -greedy method. Therefore, we showed the result of Boltzmann selection method in the result of test simulations.

V. CONCLUSION

In this paper, we proposed Q-learning algorithms in the analog model for SMC and the effectiveness of them were

TABLE IV
 THE RESULT OF OPTIMALITY OF Q-LEARNING.

		Test 1		Test 2		Test 3		Test 4		Test 5		
		No.Suc.	M.D.	No.Suc.	M.D.	No.Suc.	M.D.	No.Suc.	M.D.	No.Suc.	M.D.	
Conventional	Client	0.95	9.83	1.0	9.62	1.0	5.71	1.0	11.06	0.95	12.75	
M1	$(m = 3)$	Client	1.0	10.75	1.0	9.70	1.0	7.00	1.0	11.20	0.95	11.95
		Server $1 \sim 10$	0.05		0		0.13		0.10		0.05	
	$(m = 10)$	Client	0.95	10.31	1.0	10.17	1.0	5.66	1.0	11.24	1.0	11.47
		Server $1 \sim 10$	0		0		0.05		0		0	
M2	$(m = 3)$	Client	0.95	9.75	1.0	9.34	0.95	5.34	1.0	11.23	1.0	18.46
		Server $1 \sim 10$	0		0		0.05		0.05		0.05	
	$(m = 10)$	Client	1.0	10.32	1.0	9.83	1.0	5.29	1.0	10.72	1.0	11.46
		Server $1 \sim 10$	0		0.05		0.025		0.05		0.05	

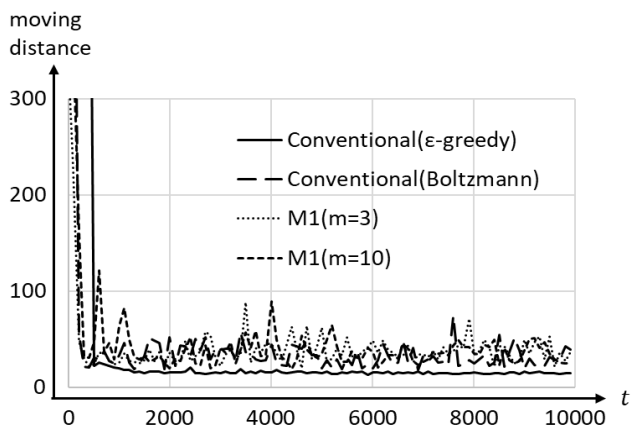


Fig. 4. The result of efficiency of Q-learning for SMC.

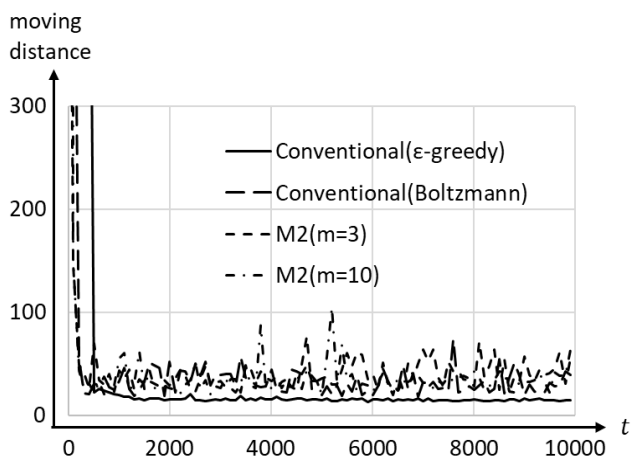


Fig. 5. The result of efficiency of Q-learning for SMC.

shown in numerical simulations. In Section 2, cloud computing system, related works on SMC and a secure data dividing mechanism used in this paper were explained. Further, Q-learning methods in the digital and analog models are introduced. In Section 3, Q-learning methods in the analog model for SMC were proposed. First, a Q-learning algorithm in the analog model was proposed using divided Q-values and the effectiveness was shown. The feature of Q-learning method for the analog model was that the action was selected as the weighted sum of plural actions. In the proposed algorithm, there was the possibility that some

servers know secure computation. Therefore, an improved Q-learning algorithm in the analog model was proposed. It was the method with dummy updating and it seems that any server is difficult to know secure computation. In section 4, numerical simulations for a maze problem were performed to show the performance of proposed methods. In the future, we will reduce the computational complexity of the client and show the safety of algorithms for SMC in theoretical proof.

REFERENCES

- [1] A. Shamir, "How to share a secret", *Comm. ACM*, Vol. 22, No. 11, pp. 612-613, 1979.
- [2] S. Subashini, and V. Kavitha, "A survey on security issues in service delivery models of cloud computing", *J. Network and Computer Applications*, Vol.34, pp.1-11, 2011.
- [3] A. Ben-David, N. Nisan, B. Pinkas, "Fair play MP: a system for secure multi-party computation", *Proceedings of the 15th ACM conference on Computer and communications security*, pp.257-266, 2008.
- [4] R. Cramer, I. Damgard, U. Maurer, "General secure multi-party computation from any linear secret-sharing scheme", *EUROCRYPT 2000*, pp.316-334, 2000.
- [5] J. Yuan, S. Yu, "Privacy Preserving Back-Propagation Neural Network Learning Made Practical with Cloud Computing", *IEEE Trans. on Parallel and Distributed Systems*, Vol.25, Issue 1, pp.212-221, 2013.
- [6] N. Rajesh, K. Sujatha, A. A. Lawrence, "Survey on Privacy Preserving Data Mining Techniques using Recent Algorithms", *International Journal of Computer Applications* Vol.133, No.7, pp.30-33, 2016.
- [7] K. Chida, et al., "A Lightweight Three-Party Secure Function Evaluation with Error Detection and Its Experimental Result", *IPSP Journal* Vol. 52 No. 9, pp. 2674-2685, Sep. 2011 (in Japanese).
- [8] H. Miyajima, N. Shigei, H. Miyajima, Y. Miyajima, S. Kitagami, N. Shiratori, "New Privacy Preserving Back Propagation Learning for Secure Multiparty Computation", *IAENG International Journal of Computer Science*, Vol.43, No.3, pp.270-276, 2016.
- [9] H. Miyajima, N. Shigei, S. Makino, H. Miyajima, Y. Miyajima, S. Kitagami, N. Shiratori, "A proposal of privacy preserving reinforcement learning for secure multiparty computation", *Artificial Intelligence Research*, Vol.6, No.2, pp.57-68, 2017.
- [10] H. Miyajima, N. Shigei, H. Miyajima, N. Shiratori, "A proposal of profit sharing method for secure multiparty computation", *International Conference on Innovative Computing, Information and Control*, 2017 (in print).
- [11] C. J. C. H. Watkins, P. Dayan, "Q-learning", *Machine Learning* 8, pp.55-68, 1992.
- [12] R. S. Sutton, A. G. Barto, "Reinforcement Learning", MIT Press, 1998.
- [13] C. Gentry, "Fully homomorphic encryption using ideal lattices", *STOC 2009*, pp.169-178, 2009.
- [14] J. Sakuma, S. Kobayashi, R. N. Wright. "Privacy-preserving reinforcement learning", In: W. W. Cohen, A. McCallum, S. T. Roweis, (eds.) *ICML*, vol. 307. ACM International Conference Proceeding Series, pp. 864-871, 2008.