

Q-learning with exploration driven by internal dynamics in chaotic neural network

1st Toshitaka Matsuki

*Department of Innovative Engineering
Faculty of Science and Technology
Oita University
Oita, Japan
matsuki@oita-u.ac.jp*

2nd Souya Inoue

*Department of Innovative Engineering
Faculty of Science and Technology
Oita University
Oita, Japan*

3rd Katsunari Shibata

*Department of Innovative Engineering
Faculty of Science and Technology
Oita University
Oita, Japan
katsunarishibata@gmail.com*

Abstract—This paper shows chaos-based reinforcement learning (RL) using a chaotic neural network (NN) functions not only with Actor-Critic, but also with Q-learning. In chaos-based RL that we have proposed, exploration is performed based on internal dynamics in a chaotic NN and the dynamics is expected to grow rational through learning. Q-learning is a very popular RL method and widely used in several researches. We focused on whether Q-learning can be adopted to chaos-based RL. Then we demonstrated the agent can learn a goal task in a grid world environment with chaos-based RL using Q-learning. It was also shown that, as learning progresses, irregularity in the network outputs originated from the internal chaotic dynamics decreases and the agent can automatically switch from exploration mode to exploitation mode. Moreover, it was confirmed that the agent can adapt to changes in the environment and automatically resume exploration.

Index Terms—Q-learning, exploration, chaotic neural network, chaos-based reinforcement learning, reservoir network

I. INTRODUCTION

Deep Learning (DL) breakthroughs have shown a deeply layered Neural Network (NN), which is trained with a large amount of data, can surpass existing image recognition systems [1] [2] [3]. DL has also been applied in several engineering fields, and many studies have identified DL as a superior approach to handcrafted systems. More recently, DL has been applied to series data processing such as speech recognition or natural language processing [4] [5] [6]. Therefore, recurrent NNs (RNN), which have a recursive structure, have been focused on. The successes of DL suggest phenomenal ability of parallel processing in the brain and difficulty in understanding and handcrafting such a massively parallel system.

Reinforcement learning (RL) is a machine learning approach in which an agent learns without any direct teacher [7]. In RL, an agent learns in a task environment by trial-and-error or exploration driven by external random noise. For a long time, End-to-End RL approach has been studied [8] [9] [10]. In End-to-End RL, an agent learns the entire process from raw input signals to action output based on RL in an NN as a parallel learning system. Mnih et al. showed that an NN can play Atari games successfully and sometimes beat human experts with RL [11]. It was also shown that an RNN can

learn tasks that require memory in which past information is integrated through time with RL [12] [13] [14]. In End-to-End RL, an NN automatically learns necessary functions as a process to generate action outputs from raw input signals through the interaction with the environment rather than given data sets. This approach is expected as a key to realize artificial general intelligence including higher functions.

In the general approach using RNN, the network requires convergence dynamics to an appropriate state to extract and memorize necessary features from the past inputs and to generate appropriate outputs. Differing from the RNN usage, the biological brain is not driven only by inputs. We consider that the state of the brain autonomously and rationally transitions as internal dynamics. We expect that organizing such dynamics in an artificial NN through learning must be a key to develop a thinking machine.

Freeman investigated chaotic dynamics that emerges in olfactory system of a rabbit [15]. In this research, it is found that olfactory bulb and cortex have many chaotic attractors that the system settles into when rabbits are held under the influence of particular odorant stimulus. When the rabbit establishes a new conditioned response to a new odorant, reorganization of existing attractors is observed and a new attractor corresponding to the new odorant is added after the emergence of chaotic dynamics. Freeman insisted that the chaos in the brain continually produces novel activity patterns that are critical to develop new nerve cell assemblies. Osana et al. demonstrated a numerical simulation of self-associative memory with a chaotic NN model that shows the similar phenomenon to the dynamics that is observed in biological experiment [16]. Freeman suggested that the ability to create activity patterns may underlie the brain's ability to generate insights and trials in trial-and-error problem solving [17].

In order to realize both convergence and transition internal dynamics in a network, which look incompatible with each other, we have investigated an approach of training an NN that has chaotic dynamics. We have set up a hypothesis that the exploration driven by disordered and autonomous state transition in a chaotic NN could develop into rational transition among multiple states through learning and the network could acquire dynamics like thinking [18]. On the hypothesis, we

This work was supported by JSPS KAKENHI Grant Number JP15K00360.

have proposed chaos-based RL in which exploration is driven by chaotic dynamics in an NN and necessary functions are acquired based on reward and punishment [19]. It has been demonstrated that a fully connected RNN having chaotic dynamics by introducing strong interconnection weights or neural refractoriness can learn an easy goal task or avoiding obstacle task with the Actor-Critic method that is a kind of RL method [20] [21]. In chaos-based RL, it was confirmed that internal chaotic dynamics initially drives exploration and the network can establish ordered dynamics to carry out a task as learning progress. Additionally, the network can automatically resume exploration when the agent meets changes in the environment. In other words, it was shown that the agent can adapt to changes in the environment and automatically switch between exploration and exploitation modes.

As well as Actor-Critic, Q-learning is a widely used popular RL method as in Deep Q-network [11]. Q-learning is a simple learning for discrete actions, and an agent learns action values for the sets of state and action. In this paper, we demonstrate that chaos-based RL can be performed with Q-learning method. Additionally, we show that an agent can automatically switch between exploration and exploitation modes with chaos based Q-learning.

Here, we employ a reservoir network (RN) that has chaotic dynamics instead of a fully connected RNN. RN was proposed as an echo state network (ESN) by Jaeger [22] or a liquid state machine by Maass [23], respectively. Its hidden layer is RNN called reservoir whose recurrent connections are randomly and sparsely connected. An RN can learn series data processing easily and stably with modifying only output layer called readout unit (RU). Hoerzer et al. showed that a chaotic RN can learn through reward-modulated Hebbian learning. This suggests that an RN can learn several tasks through exploration driven by external random noise and a sequential reward [24]. We demonstrated that an RN can learn a working memory task based on its internal chaotic dynamics and rewards, and can switch between exploration and exploitation modes automatically [25].

The remainder of this paper is organized as follows. We introduce the chaos-based RL and the experimental method in Section II. We present simulation results in Section III. Finally, we conclude these results in Section IV.

II. METHOD

A. Network

In this study, we employed an echo state network (ESN). Fig. 1 shows the network structure. This network consists of an RNN called reservoir and an output layer called readout unit. The reservoir has $N_x = 200$ neurons. The network receives $N_i = 121$ inputs as a vector $\mathbf{u}_t \in \mathbb{R}^{N_i}$, and generates $N_o = 4$ outputs as a vector $\mathbf{z}_t \in \mathbb{R}^{N_o}$. The outputs \mathbf{z}_t are fed back to the reservoir neurons. Then the internal state vector at time t is expressed as follows:

$$\mathbf{x}_t = g\mathbf{W}^{\text{rec}}\mathbf{r}_{t-1} + \mathbf{W}^{\text{in}}\mathbf{u}_t + \mathbf{W}^{\text{fb}}\mathbf{z}_{t-1}. \quad (1)$$

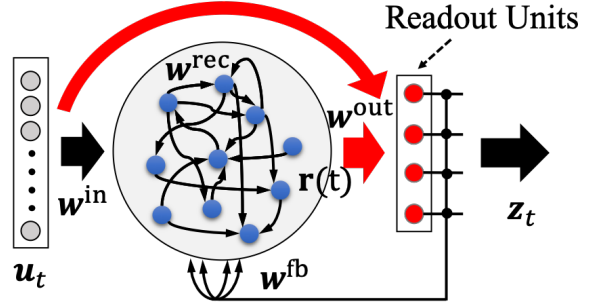


Fig. 1. Echo State Network (ESN). An ESN has a special hidden layer called “reservoir” whose neurons are randomly and sparsely connected and fixed. The output layer is called “readout units”. The inputs are fed to reservoir and readout units, and the network outputs are generated as linear combinations of reservoir activity and network inputs. The network has a feedback pathway, and outputs are fed back to the reservoir.

The initial internal state \mathbf{x}_0 is set to zero vector. $\mathbf{W}^{\text{in}} \in \mathbb{R}^{N_x \times N_i}$ is the connection weight matrix from the input to the reservoir. $\mathbf{W}^{\text{fb}} \in \mathbb{R}^{N_x \times N_o}$ is the feedback weight matrix from the network output to the reservoir. $\mathbf{W}^{\text{rec}} \in \mathbb{R}^{N_x \times N_x}$ is the recurrent connection weight matrix in the reservoir. The reservoir output vector \mathbf{r}_t is expressed as follows:

$$\mathbf{r}_t = \tanh(\mathbf{x}_t) \quad (2)$$

\mathbf{W}^{in} and \mathbf{W}^{fb} are given randomly from a uniform distribution between -1 and 1 and their connection probability is 0.5 . \mathbf{W}^{rec} is generated randomly from a uniform distribution with connection probability 0.1 and divided by its maximal absolute eigenvalue $\rho(\mathbf{W}^{\text{rec}})$ (Spectral Radius) so that it becomes a matrix with a unit spectral radius. The normalized \mathbf{W}^{rec} is scaled with $g = 1.4$, and consequently, the spectral radius is given by g . The fact is well known that the reservoir dynamics is ordered with $g < 1$, whereas it is chaotic with $g > 1$ and the larger g is, the more irregular the chaotic activity is.

The network output \mathbf{z}_t is given with the reservoir outputs \mathbf{r}_t and the input \mathbf{u}_t as follows:

$$\mathbf{z}_t = \mathbf{W}^{\text{out}} \begin{pmatrix} \mathbf{u}_t \\ \mathbf{r}_t \end{pmatrix}, \quad (3)$$

where $\mathbf{W}^{\text{out}} \in \mathbb{R}^{N_o \times (N_i + N_x)}$ is a fully connected weight matrix from the input and reservoir to the readout units, which is indicated as red arrows in Fig. 1. \mathbf{W}^{out} is initialized randomly from a uniform distribution between -0.01 and 0.01 .

In the reservoir computing approach, the readout weight \mathbf{W}^{out} is only modified to learn generating appropriate outputs. Note that, although the RN is a kind of RNN and is generally used to process series data, we use the RN as a network that has chaotic dynamics in this study. We use supervised learning to modify \mathbf{W}^{out} , and then we define the loss function as follows:

$$\mathbf{E}(\mathbf{u}_t, \mathbf{r}_t; \mathbf{W}^{\text{out}}) = \frac{1}{2}(\mathbf{T}_t - \mathbf{z}_t)^2 \quad (4)$$

where T_t is the training signal vector that is generated automatically based on Q-learning. W^{out} is modified to minimize the loss function with gradient decent as follows:

$$\begin{aligned} W^{\text{out}} &\leftarrow W^{\text{out}} - \eta \nabla_{W^{\text{out}}} E(\mathbf{u}_t, r_t; W^{\text{out}}) \\ &= W^{\text{out}} - \eta (T_t - z_t) \begin{pmatrix} \mathbf{u}_t \\ r_t \end{pmatrix}^T \end{aligned} \quad (5)$$

where $\eta = 0.01$ is the learning rate that modulates the step size of the weights update.

B. Task

Here, an agent learns to clear a simple goal task in grid world environment. The outline of the task is shown in Fig. 2. At the beginning of one episode, the agent is randomly placed on one of the four start cells in 11×11 grid world. Then, at each step, the network receives agent state input \mathbf{u}_t and generates output vector of action values $z_t = [Q_t^1, \dots, Q_t^4]^T$ whose elements are corresponding to four actions “go north”, “go east”, “go south”, “go west”, respectively. When the agent is at (i, j) in the grid world, the input vector \mathbf{u}_t is given as only one corresponding input to the agent location is 1 as follows:

$$\begin{aligned} \mathbf{u}_t &= [u_t^1, u_t^2, \dots, u_t^k, \dots, u_t^{121}]^T \\ u_t^k &= \begin{cases} 1, & k = 11(i-1) + j \\ 0, & \text{otherwise} \end{cases} \end{aligned} \quad (6)$$

Action selection is not stochastic, but always completely greedy. Then, an agent action at time t is expressed as follows:

$$a_t = \arg \max_a Q_t^a \quad (7)$$

In general RL, to find actually better actions, an agent sometimes explores with stochastic selection by taking an action that is currently estimated not to be best. For example, in ϵ -greedy policy, the agent takes random action with probability $0 \leq \epsilon \leq 1.0$. Such exploration in RL is driven by random numbers generated outside the learner agent. Meanwhile, in chaos-based RL, exploration is driven by internal chaotic dynamics in a chaotic NN. Therefore, the agent always acts with the highest action value, without any stochastic action selection.

The agent moves to the next cell in the direction according to a_t , the agent state changes from s_t to s_{t+1} and the agent receives a reward $r_{t+1}(s_t, a_t)$. This one cycle of interaction between an agent and an environment is referred to as 1 step. When an agent moves to a wall cell or outside the grid world by taking the action a_t , the agent state does not change (i.e. $s_{t+1} = s_t$) and the agent receives $r_{t+1} = -0.1$. Meanwhile, when the agent reaches the goal cell, it gets a reward $r_{t+1} = 1.0$. When the agent reaches the goal or 200 steps are passed, one episode ends and the agent is placed on one of the start cells. Then the reservoir internal state \mathbf{x}_0 is set to zero vector and new episode starts.

In this study, to examine whether the agent can adapt to changes in the environment and automatically switch between

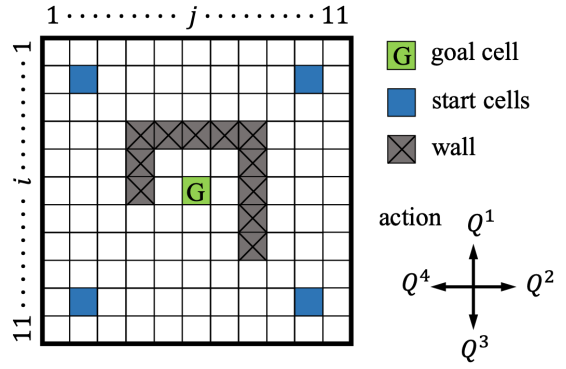


Fig. 2. The 11×11 grid world environment. Each cell corresponds to the state of the environment and can be referred by row column numbers (i, j) . At the beginning of an episode, the agent is randomly placed at the “start cell”. In one step, the agent can move one cell in four directions north, south, east and west. The agent cannot be on the “wall” and outside the world. The agent moves there and a punishment is given. If the agent reaches the “goal cell”, the episode ends with a reward.

exploration and exploitation modes, action-output correspondence is shifted on the middle of learning. More concretely, the cross arrow in Fig. 2, which shows the correspondence between action value and moving direction, rotates 90 degrees in a counter-clockwise direction.

C. Q-learning

The agent network is trained by Q-learning. When the agent takes the action a_t at state s_t and gets a reward r_{t+1} at next state s_{t+1} , the action value of a_t is evaluated as follows:

$$\hat{Q}(s_t, a_t) = r_{t+1} + \gamma \max_a Q(s_{t+1}, a) \quad (8)$$

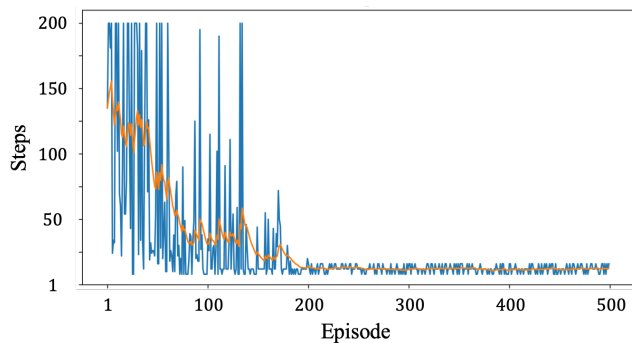
Based on this evaluated value, the action value function $Q(s_t, a_t)$ is updated as follows:

$$\begin{aligned} Q(s_t, a_t) &\leftarrow Q(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right] \\ &= (1 - \alpha)Q(s_t, a_t) + \alpha \hat{Q}(s_t, a_t) \end{aligned} \quad (9)$$

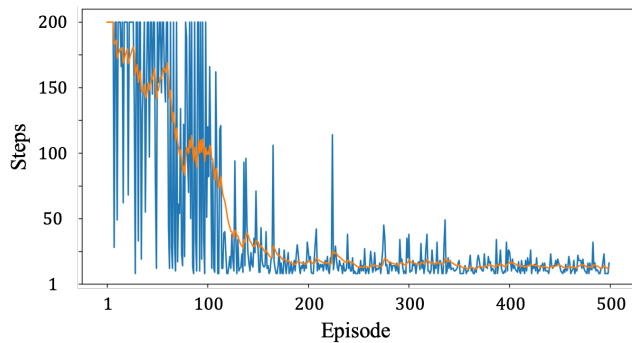
where, $0 \leq \alpha \leq 1.0$ is the step-size parameter, which influences the rate of learning. γ is a discount rate that determines the present value of discounted future rewards. In this paper, the episode ends when the agent reaches the goal, therefore, the ideal action value with which the agent reaches the goal equals to $r_{t+1} = 1.0$. Employing this algorithm to train the network in this study, the training signal for Q^{a_t} is given as follows:

$$T_t^{a_t} = r_{t+1} + \gamma \max_a Q_{t+1}^a \quad (10)$$

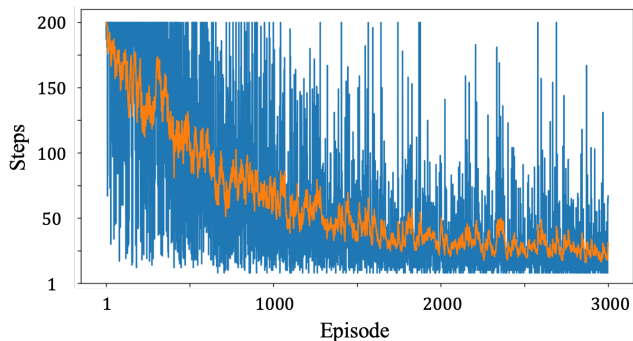
Training signals for outputs other than Q^{a_t} are generated as the errors are 0, and then, W^{out} is updated with (5). If the agent moves according to the maximum action value selected to generate the training signal of previous step in (10), the learning procedure equals to the SARSA method.



(a) Exploring with internal chaos dynamics



(b) Exploring with external random numbers

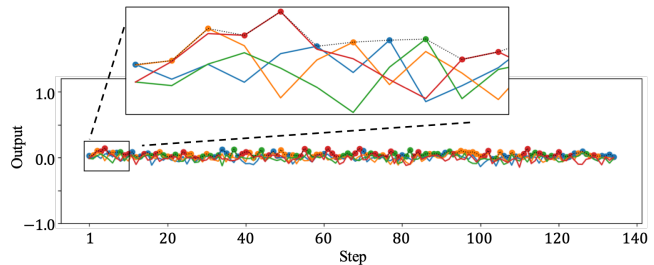


(c) The reservoir is replaced with random vector

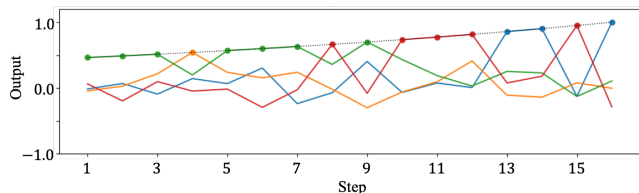
Fig. 3. Learning curve indicating the number of steps to the goal in each training episode. Orange line is a moving average of the numbers. The reduction of the number shows successful learning of the agent.

III. RESULT

Fig. 3 (a) shows the learning curve for 500 training episodes. The horizontal axis indicates the number of episodes and the vertical axis indicates the number of steps at each episode. As seen in the figure, the number of steps reduced as the learning episodes proceeded, and that shows the agent successfully learned. This result suggests that a chaotic NN can learn through Q-learning with exploration driven by its internal chaotic dynamics.



(a) episode 1



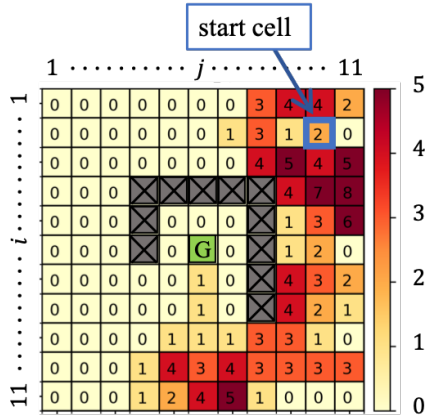
(b) episode 500

Fig. 4. Network output change in one episode for the cases of before learning (episode 1) and after learning (episode 500). Each color indicates an action as Q^1 :blue, Q^2 :orange, Q^3 :green and Q^4 :red. Plotted points and dotted line indicate the highest output at each step.

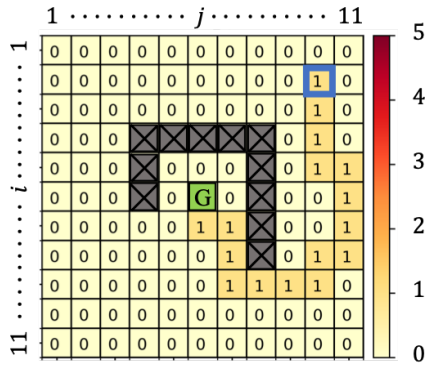
The action values that the agent network generates at the episode 1 and 500 are shown in Fig. 4 and the frequency of agent visits at the episode 1 and 500 are shown in Fig. 5. As shown in Fig. 4 (a), for early episodes, the action values chaotically fluctuate due to its chaotic internal dynamics in the network as shown in Fig. 6, and the agent seems to take actions almost randomly. Furthermore, as shown in Fig. 5 (a), the chaotic network outputs realized wide exploration of the agent in the grid world.

Meanwhile, as shown in Fig. 4 (b), the outputs became ordered and generated outputs that are almost ideal action values after learning. Fig. 5 (b) shows the route to the goal that the agent moved along. The agent successfully moved to the goal, however, we can see an unnecessary move at $(i, j) = (5, 10)$ by which the agent needs 2 extra steps than the optimal. This can be considered that the exploration was finished because of the reduction of the chaoticity before the agent had explored sufficiently in the environment.

To compare the results between random number exploration and chaos exploration, we show the training result of the same task by a feed forward NN trained with ϵ -greedy Q-learning (fixed at $\epsilon = 0.2$) in Fig. 3 (b). Here, the network has a hidden layer consisting of 200 neurons, whose activation function is tanh function, and neurons in the output layer generate outputs as linear combinations of hidden layer outputs. The network is updated with back propagation algorithm (learning rate $\eta = 0.1$) to reduce the loss function as in (4). The network succeeded to learn the task. Comparing Fig. 3 (a) and (b), a difference can be seen after reducing the number of steps until around 200 episodes. That is, in (a), the number of steps



(a) episode 1



(b) episode 500

Fig. 5. Frequency of agent visits before learning (episode 1) and after learning (episode 500). In both cases, start position is $(i, j) = (2, 10)$ (indicated by blue box).

changes by randomness of start positions in (a), meanwhile, in (b), it fluctuates more hardly than (a). It can be considered that as the learning progressed, the chaotic outputs became ordered and the exploration automatically ended in case (a) in which the exploration is driven by internal chaotic dynamics in a chaotic NN. On the other hand, in case (b), the agent continued to explore with ϵ -greedy action selection until the end of the final episode. When we adjust the ϵ value appropriately as the learning progresses, the agent stops exploring. However, the exploration in chaos based RL can be characterized by this “automatic” reduction of exploration component through learning.

Here, we wonder whether the reservoir plays the role as only a random number generator. Then, we examine the task learning with an agent network whose reservoir is replaced with N_x dimensional uniform random number vector from -1 to 1 . The learning curve is shown in Fig. 3 (c). Although the number of steps reduced as the learning progressed, it was reduced more slowly than the other cases (Fig. 3 (a) and (b)) and it had remained large and fluctuating. This result suggests that the dynamics in the reservoir is chaotic dynamics that reflects the influences of inputs from environment or learning

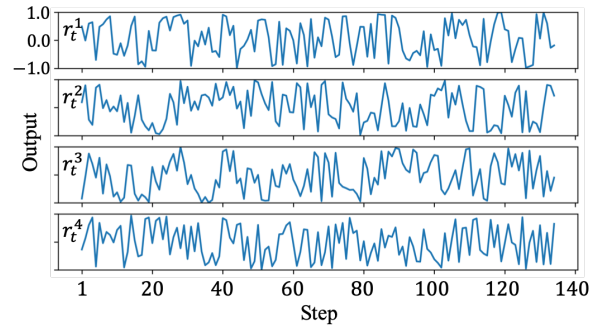
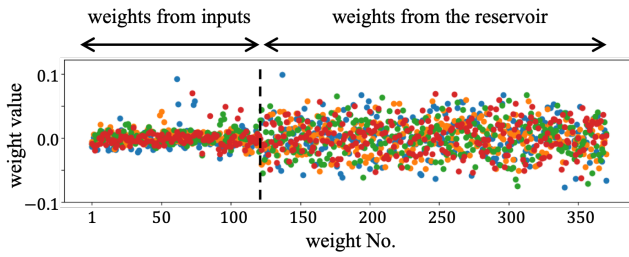


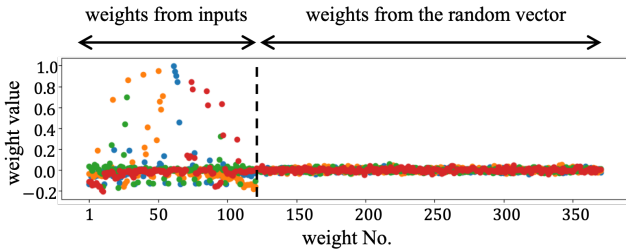
Fig. 6. Outputs of the reservoir. These lines show the values of r_t^1, \dots, r_t^4 at episode 1, respectively.

of the network rather than only random dynamics. To clarify the difference of learning between the cases (a) and (c), the training result of readout weights is shown in Fig. 7. As can be seen in Fig. 7 (a), the weights from reservoir are larger than the weights from inputs in the case of exploration driven by internal chaotic dynamics in the reservoir. This result suggests that the readout units learned to generate outputs extracting the outputs from the reservoir. Meanwhile, in the case of random vector, as shown in Fig. 7 (b), the weights for the inputs vector u_t are much larger than the weights from random vector. The reason can be considered that there is no correlation between the random numbers and the ideal outputs. These results supports the hypothesis that the reservoir is not only a generator of exploration component but also it generates the chaotic dynamics reflecting the inputs and output feedback, and we can speculate that such dynamics is useful for switching from exploration mode to exploitation mode.

We tested whether the agent can resume exploration when the rule is changed after the agent has learned under the previous rule. Here, we changed the rule every 500 episode and the result is shown in Fig. 8. The steps are rapidly increased after each rule change point because the agent could not reach the goal with the behavior that had been learned in previous episodes. However, after that, the number of steps decreased gradually. That shows that the agent successfully learned the new rules after the rule changes. Fig. 9 (a) and Fig. 10 (a) shows the network outputs and frequency of agent visits, respectively, at the episode 501, which is the first rule change point. When the agent received unknown feedback from the environment, the outputs became irregular based on the chaotic dynamics, and then the agent returns to exploration mode. Fig. 9 (b) and Fig. 10 (b) show the results at the final episode. The agent successfully learned after each rule change. This result suggests that the established connection weights, which were constructed to extract the input influence in the reservoir, were modified to learn the new rule. On the process of such reconstruction, the network dynamics became chaotic and the agent can resume exploration. These results suggest that the agent can adapt to changes in the environment and automatically switch between exploration and exploitation modes.



(a) exploration with chaos dynamics



(b) random numbers

Fig. 7. The weight values to the output layer. The horizontal axis indicates the connection weight number and the vertical axis indicates the weight values. Wight No.1 – 121 are connected from input vector u_t and the others are connected from the reservoir or the random number vector. The black dashed line is drawn at No. 121. Each plot color indicates the output neuron which the weights connect to. (a) The exploration is driven by internal dynamics in the reservoir. (b) The reservoir is replaced with a random number vector. Note that the scale of vertical axis in each figure is different.

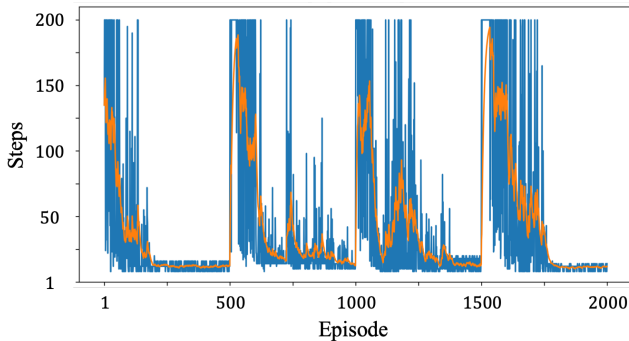
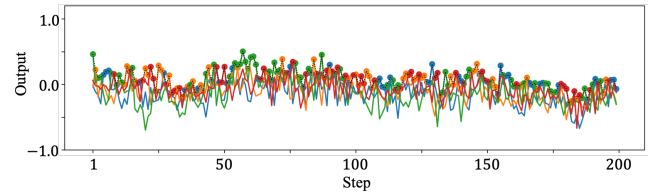


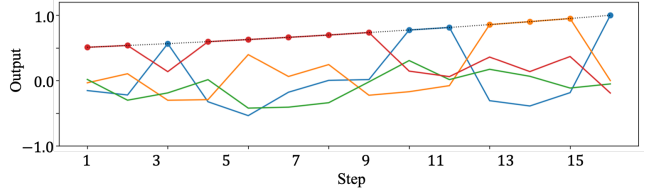
Fig. 8. Learning curve. This indicates the number of steps in each training episode. Orange line is a moving average of the numbers. The reduction of the step shows the successfully learning of the agent. The task rule is changed at episode 501, 1001 and 1501.

IV. CONCLUSION

In previous studies, we have adopted Actor-Critic method in chaos-based RL. In this paper, we demonstrated that a chaotic NN learned a goal task in a grid world environment without external random exploration noise and showed that Q-learning algorithm can be adopted in chaos-based RL. It is confirmed that irregular outputs originated from chaotic reservoir became ordered as learning progressed, and the dynamics automatically transited from exploration mode to exploitation mode. Additionally, by comparing with the result



(a) episode 501



(b) episode 2000

Fig. 9. Network output change in one episode for the cases of immediately after the first rule change (episode 501) and after learning fourth rule (episode 2000). Each color indicates an action as Q^1 :blue, Q^2 :orange, Q^3 :green and Q^4 :red. Plotted points and dotted line indicate the highest output at each step.

of the network whose reservoir is replaced with random number vector, we showed the dynamics in the reservoir reflects the inputs and output feedback, rather than only a random-like dynamics. Furthermore, it is shown that, when the task rule is changed after learning, the outputs became irregular based on the chaotic internal dynamics, and then the agent return to exploration mode and succeeded to learn the new rules.

The grid world goal task employed in this study does not require memory function for the learner agent. We predict that this frame also can be adopted to tasks that require memory function because the ESN, which is used as a network having chaotic dynamics in this study, actually has the ability to process series data. Moreover, if the network is layered more deeply, it is expected the agent can learn more complex tasks [26] [27] [28].

We consider that the phenomenon of the exploration driven by chaotic dynamics and the stabilization by learning agrees with the knowledge found by Freeman in his research about olfactory system of a rabbit. In addition, we predict that neural network dynamics around the edge of chaos is appropriate to realize transiting between chaotic and ordered states. We will examine how the network chaoticity changes as learning progresses and will investigate the relation between network chaoticity and the balance of exploration and exploitation.

REFERENCES

- [1] K. Alex, I. Sutskever, and G. E. Hinton. “Imagenet classification with deep convolutional neural networks,” *Advances in neural information processing systems*, vol. 25, pp. 1097–1105, 2012.
- [2] K. He, X. Zhang, S. Ren, J. Sun, “Deep residual learning for image recognition,” *Proc. of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [3] J. Redmon, S. Divvala, R. Girshick, A. Farhadi, “You only look once: unified, real-time object detection,” *Proc. of the IEEE conference on computer vision and pattern recognition*, pp. 779–788, 2016.

