

PAPER

A State Space Compression Method Based on Multivariate Analysis for Reinforcement Learning in High-dimensional Continuous State Spaces

Hideki SATOH^{†a)}, *Member*

SUMMARY A state space compression method based on multivariate analysis was developed and applied to reinforcement learning for high-dimensional continuous state spaces. First, useful components in the state variables of an environment are extracted and meaningless ones are removed by using multiple regression analysis. Next, the state space of the environment is compressed by using principal component analysis so that only a few principal components can express the dynamics of the environment. Then, a basis of a feature space for function approximation is constructed based on orthonormal bases of the important principal components. A feature space is thus autonomously constructed without preliminary knowledge of the environment, and the environment is effectively expressed in the feature space. An example synchronization problem for multiple logistic maps was solved using this method, demonstrating that it solves the curse of dimensionality and exhibits high performance without suffering from disturbance states.

key words: reinforcement learning, curse of dimensionality, multivariate analysis, approximation, nonlinear

1. Introduction

Reinforcement learning [1] has been receiving much attention because it can construct, without a supervisor or a model of an environment, a control function that computes an action based on the state variables of the environment. To apply reinforcement learning to the real world, we need to extend it to high-dimensional continuous state spaces because most problems in the real world comprise high-dimensional continuous states. However, as the dimension of the state space increases, the dimension of a feature space, in which a nonlinear control function is approximated, increases exponentially [1]. That is, the number of parameters to be learned increases exponentially. This is referred to as the “curse of dimensionality” and is one of the most serious problems in function approximation and reinforcement learning.

To approximate a nonlinear control function and to solve the curse of dimensionality, various approaches to autonomously constructing a feature space have been investigated. In the wire fitting approach, a feature

space for function approximation is constructed based on the control wires, that fit the control surface and create the relationship between the state and the action. Restricting the number of control wires enables the optimal action for a given state to be quickly found and reduces the number of parameters to be learned [2]. Tile coding with hashing allocates tiles whose position and size are randomly selected in a state space, and uses the tiles as a basis that constructs a feature space. By restricting the number of tiles, tile coding compresses the state space [1]. A tree-based algorithm has been developed that can handle a large continuous state space [3]. The self-organising map can perform function approximation autonomously in a feature space with a limited dimension [4]. These methods aim at effective construction of a feature space for function approximation. However, their function approximations are redundant because the basis of the feature space is not orthogonal. Because the mathematical structure of the feature space is unclear, it is difficult to ensure important properties such as the stability of the system. Hierarchical reinforcement learning is another approach to solving the curse of dimensionality. It constructs a hierarchical structure in the environment and controls each level in the hierarchy by using multiple agents. The number of parameters that each agent has to learn is reduced by dividing the state space based on the hierarchical structure [5], [6]. However, because preliminary knowledge is required in most cases to construct a hierarchical structure, it is difficult to apply hierarchical reinforcement learning to an environment whose structure is unknown. It is also difficult to mathematically clarify the demand and reason for constructing the structure.

In this paper, a method is presented that solves these problems and the curse of dimensionality. It is a state space compression method that uses multivariate analysis. It autonomously constructs a feature space without preliminary knowledge of the environment. It is based on orthogonal bases, so the mathematical structure of the system is clear in the feature space, and it is possible to mathematically analyze the dynamics and stability in a moment vector space [7]. Its ability to solve a synchronization problem for multiple logistic maps [8] using reinforcement learning demon-

Manuscript received December 22, 2005.

Manuscript revised April 4, 2006.

Final manuscript received April 17, 2006.

[†]The author is with the Future University-Hakodate, Hakodate-shi, 041-8655 Japan.

a) E-mail: jamisato@m.ieice.org

strated that the state space compression method using multivariate analysis exhibits high performance and solves the curse of dimensionality.

2. Actor-critic Method for Continuous-state Continuous-action Environment

2.1 Reinforcement Learning

A system with reinforcement learning [1] comprises two elements: an agent and an environment. The former provides a control function, and the latter denotes the target system to be controlled. The agent observes the state variables and the reward from the environment, updates the control function accordingly, and outputs a new control input to the environment. The control function and control input are referred to as policy and action, respectively.

Let us consider a discrete-time continuous-state continuous-action environment defined by a Markov decision process (MDP). Let $\mathbf{u} \stackrel{\text{def}}{=} {}^t(u_1, \dots, u_{d_u})$ be the action of dimension d_u , $\mathbf{s} \stackrel{\text{def}}{=} {}^t(s_1, \dots, s_{d_s})$ be the state of dimension d_s , $\mathcal{D}_s \stackrel{\text{def}}{=} \{s_i | s_{\min i} \leq s_i \leq s_{\min i} + D_{s_i}, 1 \leq i \leq d_s\}$, $\mathcal{D}_u \stackrel{\text{def}}{=} \{u_i | u_{\min i} \leq u_i \leq u_{\min i} + D_{u_i}, 1 \leq i \leq d_u\}$, r be the immediate reward, and superscript t denote transposition. In an MDP environment, transition probability density function $\text{prob}(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{u}_t)$ is defined to describe the relationships among \mathbf{s}_t , \mathbf{u}_t , and \mathbf{s}_{t+1} , and r_t is given when the agent takes action \mathbf{u}_t and the state changes from \mathbf{s}_t to \mathbf{s}_{t+1} . Here subscript t denotes the discrete time.

Reinforcement learning differs from supervised learning, such as neural networks, machine learning, and statistical pattern recognition. Without any help from a supervisor, it maximizes discount return R_t defined by

$$\begin{aligned} R_t &\stackrel{\text{def}}{=} \sum_{k=0}^{t_{\max}} \nu_{\text{DR}}^k r_{t+k+1} \quad \text{for } 0 < t_{\max} \leq \infty \\ &= r_{t+1} + \nu_{\text{DR}} R_{t+1}, \end{aligned} \quad (1)$$

where ν_{DR} is the discount rate ($0 \leq \nu_{\text{DR}} \leq 1$). Actor-critic methods [9], [10], which is a class of methods for solving reinforcement learning problems, consist of two parts: an actor and a critic. The actor observes state \mathbf{s}_t and decides on action \mathbf{u}_t based on policy $p(\mathbf{u}_t | \mathbf{s}_t)$, which denotes the conditional probability density function (PDF) of \mathbf{u}_t . The critic observes r_t , updates state-value function $V(\mathbf{s}) \stackrel{\text{def}}{=} E[R_t | \mathbf{s}_t = \mathbf{s}]$ accordingly, and learns policy $p(\mathbf{u} | \mathbf{s})$.

2.2 Function Approximation

Function approximation is essential for reinforcement

learning in a continuous-state continuous-action environment. It constructs an approximation of a nonlinear function. In this section, linear function approximations [1] are summarized that approximately express $V(\mathbf{s})$ and $p(\mathbf{u} | \mathbf{s})$.

State-value function $V(\mathbf{s})$ is approximated by $V(\mathbf{s}, \mathbf{c})$, which is defined by the following expansion:

$$\begin{aligned} V(\mathbf{s}, \mathbf{c}) &\stackrel{\text{def}}{=} \sum_{i=0}^{N_V} c_i \phi_{V_i}(\mathbf{s}) \\ &= {}^t \mathbf{c} \cdot \boldsymbol{\phi}_V(\mathbf{s}), \end{aligned} \quad (2)$$

where N_V denotes the degree of expansion of $V(\mathbf{s})$, $\mathbf{c} \stackrel{\text{def}}{=} {}^t(c_0, \dots, c_{N_V})$, and $\boldsymbol{\phi}_V(\mathbf{s}) \stackrel{\text{def}}{=} {}^t(\phi_{V_0}(\mathbf{s}), \dots, \phi_{V_{N_V}}(\mathbf{s}))$. Here, an orthonormal basis described in Appendix A is used for $\boldsymbol{\phi}_V(\mathbf{s})$. An $N_V + 1$ dimensional vector space with basis $\boldsymbol{\phi}_V(\mathbf{s})$ is referred to as a feature space. From Eq. (2), $V(\mathbf{s})$ is expressed by vector \mathbf{c} in the feature space with basis $\boldsymbol{\phi}_V(\mathbf{s})$.

To simplify approximation of $p(\mathbf{u} | \mathbf{s})$, let us assume that u_i obeys a normal distribution and that u_i and u_j are independent for $i \neq j$. Then, $p(\mathbf{u} | \mathbf{s})$ can be approximated by the following equation:

$$\begin{aligned} p(\mathbf{u} | \mathbf{s}) &\cong \prod_{i=1}^{d_u} p_G(u_i | \mathbf{s}, \boldsymbol{\xi}_{m_i}, \boldsymbol{\xi}_{\sigma_i}), \quad (3) \\ p_G(u_i | \mathbf{s}, \boldsymbol{\xi}_{m_i}, \boldsymbol{\xi}_{\sigma_i}) &\stackrel{\text{def}}{=} \frac{1}{(2\pi)^{1/2} \sigma_G(\mathbf{s}, \boldsymbol{\xi}_{\sigma_i})} \\ &\quad \exp\left(-\frac{(u_i - m_G(\mathbf{s}, \boldsymbol{\xi}_{m_i}))^2}{2\sigma_G(\mathbf{s}, \boldsymbol{\xi}_{\sigma_i})^2}\right), \end{aligned}$$

where $m_G(\mathbf{s}, \boldsymbol{\xi}_{m_i})$ denotes an approximation of the mean of u_i when the state of the environment is \mathbf{s} , $\sigma_G(\mathbf{s}, \boldsymbol{\xi}_{\sigma_i})$ denotes an approximation of the standard deviation of u_i when the state of environment is \mathbf{s} , $\boldsymbol{\xi}_{m_i} \stackrel{\text{def}}{=} {}^t(\xi_{m_{i0}}, \dots, \xi_{m_{iN_m}})$, $\boldsymbol{\xi}_{\sigma_i} \stackrel{\text{def}}{=} {}^t(\xi_{\sigma_{i0}}, \dots, \xi_{\sigma_{iN_\sigma}})$, N_m denotes the degree of expansion of m_G , and N_σ denotes the degree of expansion of σ_G .

Let $\boldsymbol{\phi}_m(\mathbf{s}) \stackrel{\text{def}}{=} {}^t(\phi_{m_0}(\mathbf{s}), \dots, \phi_{m_{N_m}}(\mathbf{s}))$ be an orthonormal basis. In the same manner as the approximation of $V(\mathbf{s})$, $m_G(\mathbf{s}, \boldsymbol{\xi}_{m_i})$ is defined as follows:

$$m_G(\mathbf{s}, \boldsymbol{\xi}_{m_i}) \stackrel{\text{def}}{=} {}^t \boldsymbol{\xi}_{m_i} \cdot \boldsymbol{\phi}_m(\mathbf{s}). \quad (4)$$

Because $\sigma_G(\mathbf{s}, \boldsymbol{\xi}_{\sigma_i})$ is positive, it cannot be directly expanded using an orthonormal basis, as in Eqs. (2) and (4). Thus, the following function approximation is used to transform a variable expanded by an orthonormal function expansion into a positive number:

$$\sigma_G(\mathbf{s}, \boldsymbol{\xi}_{\sigma_i}) \stackrel{\text{def}}{=} \left(\frac{\tilde{\sigma}_{\text{umax}} - \tilde{\sigma}_{\text{umin}}}{1 + \exp(-{}^t \boldsymbol{\xi}_{\sigma_i} \cdot \boldsymbol{\phi}_\sigma(\mathbf{s}))} + \tilde{\sigma}_{\text{umin}} \right) D_{u_i}, \quad (5)$$

where $\boldsymbol{\phi}_\sigma(\mathbf{s}) \stackrel{\text{def}}{=} {}^t(\phi_{\sigma_0}(\mathbf{s}), \dots, \phi_{\sigma_{N_\sigma}}(\mathbf{s}))$ is the orthonormal basis. The $\min(\sigma_G(\mathbf{s}, \boldsymbol{\xi}_{\sigma_i}))$ and $\max(\sigma_G(\mathbf{s}, \boldsymbol{\xi}_{\sigma_i}))$

are given by $\tilde{\sigma}_{\min}D_{u_i}$ and $\tilde{\sigma}_{\max}D_{u_i}$, respectively. Using these function approximations enables the mean of u_i to be expressed by vector ξ_{m_i} in a feature space with basis $\phi_m(\mathbf{s})$ and the standard deviation of u_i to be expressed by vector ξ_{σ_i} in a feature space with basis $\phi_\sigma(\mathbf{s})$.

2.3 Actor-critic Method

Let $\varepsilon_{\text{TD}t}$, \mathbf{c}_t , $\xi_{m_i;t}$, and $\xi_{\sigma_i;t}$ be the values of TD error ε_{TD} [1], \mathbf{c} , ξ_{m_i} , and ξ_{σ_i} at time t , respectively. An actor-critic method based on TD learning [1] and on the function approximations derived in Sect. 2.2 is given as follows:

Algorithm 1: Actor-critic method.

- (1-1) Initialization-1: Set \mathbf{c}_0 , $\xi_{m_i;0}$, and $\xi_{\sigma_i;0}$ for $\forall i$.
- (1-2) Initialization-2: Set $t = 0$, \mathbf{s}_0 , and \mathbf{u}_0 .
- (1-3) Environment: Receive \mathbf{u}_t , and output r_{t+1} and \mathbf{s}_{t+1} to the agent.
- (1-4) Actor-critic:
 - (1) Critic-1: Calculate $\varepsilon_{\text{TD}t}$ and \mathbf{c}_{t+1} by using Eqs. (9) and (11) (Update $V(\mathbf{s}, \mathbf{c})$).
 - (2) Critic-2: Calculate $\xi_{m_i;t+1}$, $\xi_{\sigma_i;t+1}$ by using Eqs. (12) and (13) (Update $p_G(u_i|\mathbf{s}, \xi_{m_i}, \xi_{\sigma_i})$) for $\forall i$.
 - (3) Actor: Set $u_{i;t+1}$ based on $p_G(u_{i;t+1}|\mathbf{s}_{t+1}, \xi_{m_i;t+1}, \xi_{\sigma_i;t+1})$ for $\forall i$.
- (1-5) Set $t = t + 1$. If an episode finished, go to Step (1-2), else go to Step (1-3).

Here, episode denotes a trial [1]. In Step (1-1), initial values \mathbf{c}_0 , $\xi_{m_i;0}$, and $\xi_{\sigma_i;0}$ are set for $\forall i$ as follows:

$$\begin{aligned} \mathbf{c}_0 &= \mathbf{0}, \\ \xi_{m_{ij};0} &= \begin{cases} u_{\min i} + D_{u_i}/2 & \text{for } j = 0 \\ 0 & \text{for } 1 \leq j \leq N_m, \end{cases} \\ \xi_{\sigma_{ij};0} &= \begin{cases} \xi_0 & \text{for } j = 0 \\ 0 & \text{for } 1 \leq j \leq N_\sigma, \end{cases} \end{aligned}$$

where ξ_0 is set so that $\sigma_G(\mathbf{s}, \xi_{\sigma_i})$ takes an appropriate value for initial learning. In Sect. 5, ξ_0 is set to -1.0 for $\forall i$, so that $\sigma_G(\mathbf{s}_0, \xi_{\sigma_i;0}) = 0.135D_{u_i}$.

Let $p(\mathbf{s})$ be the PDF of \mathbf{s} and $MSE(\mathbf{c})$ be the mean square error (MSE) of state-value function $V(\mathbf{s})$ defined by

$$MSE(\mathbf{c}) \stackrel{\text{def}}{=} \int_{\mathbf{s} \in \mathcal{D}_s} p(\mathbf{s})(V(\mathbf{s}) - V(\mathbf{s}, \mathbf{c}))^2 d\mathbf{s}. \quad (6)$$

A method [1] that minimizes $MSE(\mathbf{c})$ with respect to \mathbf{c} and updates $V(\mathbf{s}, \mathbf{c})$ is summarized below. Let $\hat{V}(\mathbf{s})$ be an estimation of the state-value function. We can obtain $\hat{V}(\mathbf{s})$ by using the following equation:

$$\hat{V}(\mathbf{s}_t) = r_{t+1} + \nu_{\text{DR}} V(\mathbf{s}_{t+1}, \mathbf{c}_t). \quad (7)$$

Thus, by replacing \mathbf{x} with \mathbf{s} , ζ with ζ_V , g with V , and \hat{g} with \hat{V} in Eq. (A.8), we obtain the following equation:

$$\mathbf{c}_{t+1} = \mathbf{c}_t + \zeta_V (\hat{V}(\mathbf{s}_t) - V(\mathbf{s}_t, \mathbf{c}_t)) \nabla_c V(\mathbf{s}_t, \mathbf{c}_t), \quad (8)$$

where ζ_V is the step size. By substituting TD error $\varepsilon_{\text{TD}t}$, defined by

$$\begin{aligned} \varepsilon_{\text{TD}t} &\stackrel{\text{def}}{=} \hat{V}(\mathbf{s}_t) - V(\mathbf{s}_t, \mathbf{c}_t) \\ &= r_{t+1} + \nu_{\text{DR}} V(\mathbf{s}_{t+1}, \mathbf{c}_t) - V(\mathbf{s}_t, \mathbf{c}_t), \end{aligned} \quad (9)$$

and $\nabla_c V(\mathbf{s}, \mathbf{c}) = \phi_V(\mathbf{s})$ into Eq. (8), we obtain

$$\mathbf{c}_{t+1} = \mathbf{c}_t + \zeta_V \varepsilon_{\text{TD}t} \phi_V(\mathbf{s}_t). \quad (10)$$

Let ν_{Tr} be the trace-decay parameter [1]. By introducing eligibility trace $\boldsymbol{\eta}_t$ [1] into the above equation, we obtain the following equation for updating \mathbf{c} :

$$\begin{aligned} \boldsymbol{\eta}_t &= \nu_{\text{DR}} \nu_{\text{Tr}} \boldsymbol{\eta}_{t-1} + \phi_V(\mathbf{s}_t), \\ \mathbf{c}_{t+1} &= \mathbf{c}_t + \zeta_V \varepsilon_{\text{TD}t} \boldsymbol{\eta}_t. \end{aligned} \quad (11)$$

Here, $\boldsymbol{\eta}_0$, the initial value of $\boldsymbol{\eta}_t$, is set to $\mathbf{0}$.

In the same manner as for $V(\mathbf{s}, \mathbf{c})$, $p_G(u_i|\mathbf{s}, \xi_{m_i}, \xi_{\sigma_i})$ can be updated. Let ζ_m be the step size to update $\xi_{m_i;t}$, and ζ_σ be the step size to update $\xi_{\sigma_i;t}$. By replacing \mathbf{c}_t with $\xi_{m_i;t}$, \tilde{g} with $u_{i;t}$, $g(\mathbf{x}_t, \mathbf{c}_t)$ with $m_G(\mathbf{s}_t, \xi_{m_i;t})$, and $\nabla_c g(\mathbf{x}, \mathbf{c})$ with $\nabla_{\xi_{m_i}} \xi_{m_i} \cdot \phi_m(\mathbf{s}_t) = \phi_m(\mathbf{s}_t)$ in Eq. (A.9), we obtain the following equation:

$$\begin{aligned} \xi_{m_i;t+1} &= \xi_{m_i;t} \\ &\quad + \zeta_m \varepsilon_{\text{TD}t} (u_{i;t} - m_G(\mathbf{s}_t, \xi_{m_i;t})) \phi_m(\mathbf{s}_t). \end{aligned} \quad (12)$$

In the same manner as for Eq. (12), by replacing \mathbf{c}_t with $\xi_{\sigma_i;t}$, \tilde{g} with $|u_{i;t} - m_G(\mathbf{s}_t, \xi_{m_i;t})|$, $g(\mathbf{x}_t, \mathbf{c}_t)$ with $\sigma_G(\mathbf{s}_t, \xi_{\sigma_i;t})$, and $\nabla_c g(\mathbf{x}, \mathbf{c})$ with $\nabla_{\xi_{\sigma_i}} \xi_{\sigma_i} \cdot \phi_\sigma(\mathbf{s}_t) = \phi_\sigma(\mathbf{s}_t)$ in Eq. (A.9), we obtain the following equation:

$$\begin{aligned} \xi_{\sigma_i;t+1} &= \xi_{\sigma_i;t} + \zeta_\sigma \varepsilon_{\text{TD}t} (|u_{i;t} - m_G(\mathbf{s}_t, \xi_{m_i;t})| \\ &\quad - \sigma_G(\mathbf{s}_t, \xi_{\sigma_i;t})) \phi_\sigma(\mathbf{s}_t). \end{aligned} \quad (13)$$

3. State Space Compression Method and its Application to Reinforcement Learning

3.1 Framework of State Space Compression and Learning

Let us consider a system in which several controllers control an environment, and assume that the agent is one of the controllers. The controllers observe the state variables and the reward from the environment, and output control inputs to the environment. Suppose that the agent can observe state \mathbf{x}_t , reward r_{t+1} , and control input \mathbf{v}_t that includes control inputs determined by the other controllers and action \mathbf{u}_t determined by the agent. Note that it is assumed that \mathbf{x}_t and \mathbf{v}_t does not not always include all the states in the environment and all the control inputs determined by the

other controllers. To simplify the equations described later, the relationship between \mathbf{u}_t and \mathbf{v}_t is defined as follows:

$$\mathbf{u}_t = G_u \mathbf{v}_t + \mathbf{g}_u. \quad (14)$$

Here, $\mathbf{x}_t \stackrel{\text{def}}{=} {}^t(x_{1;t}, \dots, x_{d_x;t})$, $\mathbf{v}_t \stackrel{\text{def}}{=} {}^t(v_{1;t}, \dots, v_{d_v;t})$, $\mathbf{D}_x \stackrel{\text{def}}{=} \{x_i \mid x_{\min i} \leq x_i \leq x_{\min i} + D_{x_i}, 1 \leq i \leq d_x\}$, $\mathbf{D}_v \stackrel{\text{def}}{=} \{v_i \mid v_{\min i} \leq v_i \leq v_{\min i} + D_{v_i}, 1 \leq i \leq d_v\}$, and G_u and \mathbf{g}_u are a matrix and a vector, respectively, obtained from $v_{\min i}$, D_{v_i} , $u_{\min i}$, and D_{u_i} .

When the basis of \mathbf{x} for function approximation is constructed using the direct product of the bases of x_i , the degree of expansion of \mathbf{x} increases proportionally to the d_x th power of the degree of expansion of x_i (see Eq. (A.4)). This is referred to as the ‘‘curse of dimensionality.’’ The framework of state space compression and reinforcement learning presented in this section solves this problem. In the framework, the agent compresses the state space, constructs a feature space with a limited dimension, extracts the principal dynamics in \mathbf{x}_t , and determines \mathbf{u}_t . The procedure operated by the agent is summarized as follows (the numbers in parentheses correspond to (1), \dots , (7) in Step (2-4) of Algorithm 2). (3) Compute statistics such as mean and covariance of \mathbf{x} . (4) Compute the resolution (which denotes the importance) of x_i using multiple regression analysis. Based on the resolution, extract useful components of \mathbf{x} and remove meaningless ones. (5) Compute principal axes matrix M , which maps the useful components of \mathbf{x} to principal component vector $\mathbf{y} = {}^t(y_1, \dots, y_{d_y})$, which express the principal dynamics in the useful components of \mathbf{x} . (6) Compute the resolution of y_i . (7) Based on the resolution of y_i , determine bases ϕ_v , ϕ_m , and ϕ_σ . (1) Transform \mathbf{x} into \mathbf{y} using principal axes matrix M , and normalize \mathbf{y} into normalized principal state \mathbf{s} so that \mathbf{s} contains the principal dynamics in \mathbf{x} in a bounded space. (2) Compute action \mathbf{u} based on \mathbf{s} using the actor-critic method in the feature spaces with bases ϕ_v , ϕ_m , and ϕ_σ .

In the above procedure, it is necessary, before performing the actor-critic method, to obtain principal axes matrix M , bases ϕ_v , ϕ_m , and ϕ_σ , and normalized principal state \mathbf{s} . Thus, before performing the actor-critic method, \mathbf{u} is determined at random, and statistics such as the mean and covariance of \mathbf{x} are computed. After a sufficiently long time has elapsed, M , ϕ_v , ϕ_m , and ϕ_σ are computed. They are updated at appropriate intervals so that they include the latest dynamics in the environment. The time when M , ϕ_v , ϕ_m , and ϕ_σ are computed is indicated by using two variables: n_{MA} and trg_{MA} . The former denotes the number of times M , ϕ_v , ϕ_m , and ϕ_σ have been updated, and the latter denotes the trigger to update M , ϕ_v , ϕ_m , and ϕ_σ . Here, M , ϕ_v , ϕ_m , and ϕ_σ are updated at regular interval T_{RB} . The structure and algorithm for the agent are shown in Fig. 1 and given by Algorithm 2, respectively.

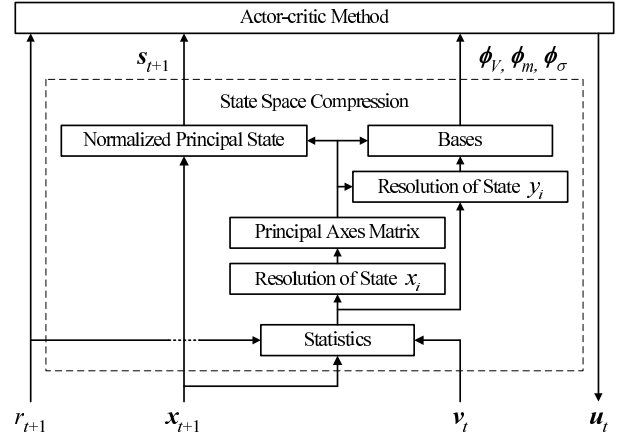


Fig. 1 Structure of agent with state space compression.

Algorithm 2: Reinforcement learning with state space compression.

- (2-1) Initialization-1: Set $\text{trg}_{\text{MA}} = 0$, $n_{\text{MA}} = 0$, \mathbf{c}_0 , $\xi_{\mathbf{m}0}$, and $\xi_{\sigma 0}$.
- (2-2) Initialization-2: Set $t = 0$, \mathbf{x}_0 , and \mathbf{u}_0 .
- (2-3) Environment: Receive \mathbf{u}_t , and output r_{t+1} and \mathbf{x}_{t+1} to the agent.
- (2-4) Agent:
 - (1) Transform \mathbf{x}_{t+1} into \mathbf{s}_{t+1} by using Eqs. (24) and (32) if $n_{\text{MA}} \geq 1$.
 - (2) Perform actor-critic method (Step (1-4) in Algorithm 1) if $n_{\text{MA}} \geq 1$, or set \mathbf{u}_t randomly if $n_{\text{MA}} = 0$.
 - (3) Compute statistics of \mathbf{x} , and set n_{MA} and trg_{MA} by using Algorithm 3.
 - (4) Compute resolution of state x_i by using Eq. (20) if $\text{trg}_{\text{MA}} = 1$.
 - (5) Update principal axes matrix and mean by using Eqs. (23) and (25) if $\text{trg}_{\text{MA}} = 1$.
 - (6) Compute resolution of y_i by using Eq. (31) if $\text{trg}_{\text{MA}} = 1$.
 - (7) Update bases as shown in Sect. 3.7 and set $\text{trg}_{\text{MA}} = 0$ if $\text{trg}_{\text{MA}} = 1$.
- (2-5) Set $t = t + 1$. If an episode finished, go to Step (2-2), else go to Step (2-3).

Each step of the agent is described in detail below.

3.2 Statistics of Environment

To enable multivariate analysis to be performed, the mean and covariance of \mathbf{x}_t , \mathbf{x}_{t+1} , \mathbf{v}_t , and r_{t+1} are computed at appropriate intervals by using Algorithm 3 so that they include the latest dynamics in the environment, where COV_{xx} denotes the covariance matrix of

$(x_{1;\tau}, \dots, x_{d_x;\tau})$, COV_{vxx} denotes the covariance matrix of $(v_{1;\tau}, \dots, v_{d_v;\tau}, x_{1;\tau}, \dots, x_{d_x;\tau})$, $COV_{vxx'vxx'}$ denotes the covariance matrix of $(v_{1;\tau}, \dots, v_{d_v;\tau}, x_{1;\tau}, \dots, x_{d_x;\tau}, x_{1;\tau+1}, \dots, x_{d_x;\tau+1})$, $\mathbf{cov}_{x'vxi} \stackrel{\text{def}}{=} {}^t(\text{cov}(x_{i;\tau+1}, v_{1;\tau}), \dots, \text{cov}(x_{i;\tau+1}, v_{d_v;\tau}), \text{cov}(x_{i;\tau+1}, x_{1;\tau}), \dots, \text{cov}(x_{i;\tau+1}, x_{d_x;\tau}))$, $\mathbf{cov}_{rvxx'} \stackrel{\text{def}}{=} {}^t(\text{cov}(r_{\tau+1}, v_{1;\tau}), \dots, \text{cov}(r_{\tau+1}, v_{d_v;\tau}), \text{cov}(r_{\tau+1}, x_{1;\tau}), \dots, \text{cov}(r_{\tau+1}, x_{d_x;\tau}), \text{cov}(r_{\tau+1}, x_{1;\tau+1}), \dots, \text{cov}(r_{\tau+1}, x_{d_x;\tau+1}), \text{cov}(x_{\tau}, y_{\tau}))$ denotes the covariance of x_{τ} and y_{τ} , $\bar{\mathbf{x}}$ denotes the mean of \mathbf{x}_{τ} , and $t - T_S + 1 \leq \tau \leq t$.

Algorithm 3: Compute Statistics.

- (3-1) Initialization: If $n_{MA} = 0$, set $T = T_S$. If $n_{MA} \geq 0$, set $T = T_{RB}$.
- (3-2) If the remainder of $t + 1$ upon division by T is equal to 0 (i.e., at periodic intervals of T), perform the following steps.
- (3-3) Compute COV_{xx} , COV_{vxx} , $COV_{vxx'vxx'}$, $\mathbf{cov}_{x'vxi}$, $\mathbf{cov}_{rvxx'}$, and $\bar{\mathbf{x}}$.
- (3-4) Set $n_{MA} = n_{MA} + 1$ and $trg_{MA} = 1$.

Here, the statistics of \mathbf{x} are initially computed at $t = T_S - 1$. Thereafter, the bases are updated at periodic intervals of T_{RB} .

3.3 Resolution of State x_i

To determine action \mathbf{u}_t , it is necessary to estimate which elements in state vector \mathbf{x}_t are important. The resolution of $x_{i;t}$ for computing \mathbf{u}_t is estimated by using the partial regression coefficients [11] and standard deviations.

First, the effect of \mathbf{v}_t and \mathbf{x}_t on \mathbf{x}_{t+1} is estimated by using the partial regression coefficients for \mathbf{x}_{t+1} . The multiple regression equation [11] for criterion variable \mathbf{x}_{t+1} with respect to explanatory variables \mathbf{v}_t and \mathbf{x}_t is expressed by

$$\mathbf{x}_{t+1} = \beta_0 + B_v \mathbf{v}_t + B_x \mathbf{x}_t, \quad (15)$$

where β_{i0} , β_{vij} , and β_{xij} are the partial regression coefficients, $\beta_0 \stackrel{\text{def}}{=} {}^t(\beta_{10}, \dots, \beta_{d_x0})$, $\beta_{vi} \stackrel{\text{def}}{=} {}^t(\beta_{vi1}, \dots, \beta_{vid_v})$, $\beta_{xi} \stackrel{\text{def}}{=} {}^t(\beta_{xi1}, \dots, \beta_{xid_x})$, $B_v \stackrel{\text{def}}{=} {}^t[\beta_{v1}, \dots, \beta_{vd_x}]$ is a $d_x \times d_v$ matrix, and $B_x \stackrel{\text{def}}{=} {}^t[\beta_{x1}, \dots, \beta_{xd_x}]$ is a $d_x \times d_x$ matrix. Partial regression coefficients β_{vij} and β_{xij} are derived by the following equation [11]:

$$\beta_{vxi} = COV_{vxx}^{-1} \mathbf{cov}_{x'vxi}, \quad (16)$$

where $\beta_{vxi} \stackrel{\text{def}}{=} {}^t({}^t\beta_{vi}, {}^t\beta_{xi})$.

Next, the effect of \mathbf{v}_t , \mathbf{x}_t , and \mathbf{x}_{t+1} on r_{t+1} is estimated by using the partial regression coefficient for r_{t+1} . The multiple regression equation for criterion variable r_{t+1} with respect to explanatory variables \mathbf{v}_t , \mathbf{x}_t , and \mathbf{x}_{t+1} is expressed by

$$r_{t+1} = b_0 + {}^t\mathbf{b}_v \mathbf{v}_t + {}^t\mathbf{b}_x \mathbf{x}_t + {}^t\mathbf{b}_{x'} \mathbf{x}_{t+1}, \quad (17)$$

where b_0 , b_{vj} , b_{xj} , and $b_{x'j}$ are the partial regression coefficients, $\mathbf{b}_v \stackrel{\text{def}}{=} {}^t(b_{v1}, \dots, b_{vd_v})$, $\mathbf{b}_x \stackrel{\text{def}}{=} {}^t(b_{x1}, \dots, b_{xd_x})$, and $\mathbf{b}_{x'} \stackrel{\text{def}}{=} {}^t(b_{x'1}, \dots, b_{x'd_x})$. In the same manner as in Eq. (16), partial regression coefficients b_{vj} , b_{xi} , and $b_{x'i}$ are derived by the following equation:

$$\mathbf{b}_{vxx'} = (COV_{vxx'vxx'})^{-1} \mathbf{cov}_{rvxx'}, \quad (18)$$

where $\mathbf{b}_{vxx'} \stackrel{\text{def}}{=} {}^t({}^t\mathbf{b}_v, {}^t\mathbf{b}_x, {}^t\mathbf{b}_{x'})$.

Let res_{xvij} be the resolution of x_i for computing v_j . A combination of $\sigma(x_i)$, D_{vj} , $|\beta_{vij}|$, and $|b_{xi}| + |b_{x'i}|$ is used to determine res_{xvij} for three reasons:

- (1) The effect of $x_{i;t}$ on $v_{j;t}$ denotes the agent itself. Thus, if the effect of $x_{i;t}$ on $v_{j;t}$ is used to determine the resolution, a positive feedback loop arises, and a system with the agent may become unstable. Thus, instead of the effect of $x_{i;t}$ on $v_{j;t}$, $|\beta_{vij}|$, which denotes the effect of $v_{j;t}$ on $x_{i;t+1}$, is used.
- (2) Important information about the structure and dynamics of the system are contained in $\sigma(x_i)$ and $\sigma(v_j)$. However, $\sigma(v_j)$ contains the motion in the agent, so the system may become unstable if $\sigma(v_j)$ is used for the same reason described above. Thus, D_{vj} instead of $\sigma(v_j)$ and $\sigma(x_i)$ are used.
- (3) Because the agent works to obtain a large reward, it is effective to introduce $|b_{xi}| + |b_{x'i}|$, which denotes the effect of $x_{i;t}$ and $x_{i;t+1}$ on r_{t+1} .

Based on these reasons, res_{xvij} is determined by

$$res_{xvij} = \sigma(x_i)^{\gamma_1} D_{vj}^{\gamma_2} |\beta_{vij}|^{\gamma_3} (|b_{xi}| + |b_{x'i}|)^{\gamma_4}, \quad (19)$$

where $\gamma_1, \dots, \gamma_4$ are given constants used to adjust the combination. Here, $\sigma(x_i) = \text{cov}(x_i, x_i)^{1/2}$, and $\sigma(x_i)$ is obtained in Step (3-3) of Algorithm 3.

Let res_{xui} be the resolution of x_i used for computing \mathbf{u}_t . As shown in Sect. 3.1, the relationship between \mathbf{u}_t and \mathbf{v}_t is defined by Eq. (14). Thus, res_{xui} is determined by the weighted sum of res_{xvij} with respect to j as follows:

$$res_{xui} = \sum_{j=1}^{d_u} (G_u {}^t(res_{xvi1}, \dots, res_{xvid_x}))_j, \quad (20)$$

where $(\cdot)_j$ denotes entry j of a vector.

3.4 Principal Axes Matrix

As mentioned in Sect. 3.1, when the dimension of \mathbf{x}_t is high, the degree of expansion of \mathbf{x}_t is too large to perform reinforcement learning. Thus, it is necessary to express the dynamics in \mathbf{x}_t in a low-dimensional feature space. In this section, principal components, which represent the important dynamics, are derived. Let $y_{i;t}$ be the i th principal component of \mathbf{x}_t , let \mathbf{y}_t be ${}^t(y_{1;t}, \dots, y_{d_y;t})$. The importance of y_i is computed in

Sect. 3.5 and y_i used for control is determined in Sect. 3.7. Thus, d_y is set to d_x .

Principal component vector \mathbf{y}_t is obtained as the product of principal axes matrix M and state vector \mathbf{x}_t , and M is derived from the eigenvector of COV_{xx} [11]. Most of the dynamics in \mathbf{x}_t can be expressed by using only a small number of elements in \mathbf{y}_t if $x_{i;t}$ and $x_{j;t}$ for $i \neq j$ are correlated. However, it is not always true that elements in \mathbf{x}_t with large dynamics contribute to reward r_{t+1} . Some elements may not be related to r_{t+1} . To remove the effect of such elements on \mathbf{y}_t , the covariance matrix for $res_{xu_i}x_i$ is used instead of that for x_i . Let $\mathbf{x}^R_t \stackrel{\text{def}}{=} {}^t(res_{xu_1}x_{1;t}, \dots, res_{xu_{d_x}}x_{d_x;t})$ and COV_{xx}^R be the covariance matrix of \mathbf{x}^R . From the definition of the covariance, COV_{xx}^R is obtained by using the elements of COV_{xx} as follows:

$$COV_{xx}^R = [res_{xu_i}res_{xu_j}cov(x_i, x_j)]. \quad (21)$$

Let \mathbf{e}^R_i be the i th eigen vector of COV_{xx}^R and $M^R \stackrel{\text{def}}{=} [\mathbf{e}^R_1, \dots, \mathbf{e}^R_{d_x}]$ be the principal axes matrix of \mathbf{x}^R_t . We then obtain principal component vector \mathbf{y}_t as follows:

$$\begin{aligned} \mathbf{y}_t &= {}^tM^R \mathbf{x}^R_t \\ &= {}^tM^R \text{diag}[res_{xu_i}] \mathbf{x}_t, \end{aligned} \quad (22)$$

where $\text{diag}[res_{xu_i}]$ denotes the diagonal matrix of res_{xu_i} . Let principal axes matrix M of \mathbf{x}_t be defined by the following equation:

$$M \stackrel{\text{def}}{=} \text{diag}[res_{xu_i}] M^R. \quad (23)$$

By substituting Eq. (23) into Eq. (22), we obtain

$$\mathbf{y}_t = {}^tM \mathbf{x}_t. \quad (24)$$

Let λ_i be the variance in $y_{i;t}$, and $\bar{\mathbf{y}}$ be the mean of \mathbf{y}_t . The former is equal to the i th eigen value of COV_{xx}^R , and the latter is derived by

$$\bar{\mathbf{y}} = {}^tM \bar{\mathbf{x}}. \quad (25)$$

3.5 Resolution of State y_i

As shown in Appendix A, bases ϕ_v , ϕ_m , and ϕ_σ are determined by the degree of expansion of y_i . Thus, to obtain a large discount return, it is reasonable to assign a larger value to the degree of expansion of $y_{i;t}$ if $y_{i;t}$ is closely related to \mathbf{u}_t and r_{t+1} . The resolution of $y_{i;t}$ is derived so as to achieve this aim in the same manner as in Sect. 3.3.

The multiple regression equation for \mathbf{y}_{t+1} is described in the same manner as in Eq. (15) as follows:

$$\mathbf{y}_{t+1} = \tilde{\beta}_0 + \tilde{B}_v \mathbf{v}_t + \tilde{B}_y \mathbf{y}_t, \quad (26)$$

where $\tilde{\beta}_{i0}$, $\tilde{\beta}_{vij}$, and $\tilde{\beta}_{yij}$ are the partial regression coefficients, $\tilde{\beta}_0 \stackrel{\text{def}}{=} {}^t(\tilde{\beta}_{10}, \dots, \tilde{\beta}_{d_x0})$, $\tilde{\beta}_{vi} \stackrel{\text{def}}{=} {}^t(\tilde{\beta}_{vi1}, \dots, \tilde{\beta}_{vid_v})$, $\tilde{\beta}_{yi} \stackrel{\text{def}}{=} {}^t(\tilde{\beta}_{yi1}, \dots, \tilde{\beta}_{yid_y})$, $\tilde{B}_v \stackrel{\text{def}}{=} {}^t[\tilde{\beta}_{v1}, \dots, \tilde{\beta}_{vd_y}]$ is a $d_y \times d_v$ matrix, and $\tilde{B}_y \stackrel{\text{def}}{=} {}^t[\tilde{\beta}_{y1}, \dots, \tilde{\beta}_{yd_y}]$ is a $d_y \times d_y$ matrix. By substituting Eq. (24) into Eq. (15), we obtain

$$\mathbf{y}_{t+1} = {}^tM \beta_0 + {}^tM B_v \mathbf{v}_t + {}^tM B_x M \mathbf{y}_t. \quad (27)$$

By comparing Eqs. (26) and (27), partial regression coefficients $\tilde{\beta}_{vij}$ is obtained by the following equation:

$$\tilde{B}_v = {}^tM B_v. \quad (28)$$

By substituting Eq. (24) into Eq. (17), we can describe the multiple regression equation for r_{t+1} by using

$$\begin{aligned} r_{t+1} &= b_0 + \mathbf{b}_v \mathbf{v}_t + {}^t\mathbf{b}_x M \mathbf{y}_t + {}^t\mathbf{b}_{x'} M \mathbf{y}_{t+1} \\ &= b_0 + \mathbf{b}_v \mathbf{v}_t + {}^t\tilde{\mathbf{b}}_y \mathbf{y}_t + {}^t\tilde{\mathbf{b}}_{y'} \mathbf{y}_{t+1}. \end{aligned} \quad (29)$$

Here, ${}^t\tilde{\mathbf{b}}_y \stackrel{\text{def}}{=} ({}^t\tilde{b}_{y1}, \dots, {}^t\tilde{b}_{yd_y}) = {}^t\mathbf{b}_x M$ and ${}^t\tilde{\mathbf{b}}_{y'} \stackrel{\text{def}}{=} ({}^t\tilde{b}_{y'1}, \dots, {}^t\tilde{b}_{y'd_y}) = {}^t\mathbf{b}_{x'} M$. As a result, we obtain partial regression coefficients \tilde{b}_{y_i} and $\tilde{b}_{y'_i}$ from

$$\begin{aligned} \tilde{b}_y &= {}^tM \mathbf{b}_x, \\ \tilde{b}_{y'} &= {}^tM \mathbf{b}_{x'}. \end{aligned} \quad (30)$$

Let $res_{yv_{ij}}$ be the resolution of y_i for computing v_j and res_{yu_i} be the resolution of y_i for computing \mathbf{u} . In the same manner as in Eqs. (19) and (20), res_{yu_i} is determined as follows:

$$res_{yu_i} = \sum_{j=1}^{d_u} (G_u(res_{yv_{i1}}, \dots, res_{yv_{id_x}}))_j, \quad (31)$$

$$res_{yv_{ij}} = \sigma(y_i)^{\tilde{\gamma}_1} D_{v_j} \tilde{\gamma}_2 |\tilde{\beta}_{vij}|^{\tilde{\gamma}_3} (|\tilde{b}_{y_i}| + |\tilde{b}_{y'_i}|)^{\tilde{\gamma}_4},$$

where $\tilde{\gamma}_1, \dots, \tilde{\gamma}_4$ are given constants, and $\sigma(y_i) = \lambda_i^{1/2}$.

3.6 State Transformation

Principal component vector \mathbf{y}_t obtained by using Eq. (24) has a unique boundary because it is assumed that $\mathbf{x}_t \in \mathbf{D}_x$, as mentioned in Sect. 3.1. However, \mathbf{y}_t does not take all the values within the boundary. Thus, state vector \mathbf{s}_t used by the actor-critic method is set by using the following equation so that it is distributed over the whole feature space and expresses the dynamics in \mathbf{x}_t as accurately as possible:

$$s_{i;t} = \frac{D_{s_i}}{1 + \exp(-\frac{y_{i;t} - \bar{y}_i}{c_{PN} \sigma(y_i)})} - \frac{D_{s_i}}{2} + \bar{s}_i, \quad (32)$$

where $d_s = d_y$, $\bar{s}_i = s_{\min i} + D_{s_i}/2$, and c_{PN} is set to 0.6 so that the first term on the right side of Eq. (32) is approximately equal to the probability distribution function of a normal distribution with mean \bar{y}_i and standard deviation $\sigma(y_i)$.

3.7 Basis of Feature Space

Based on the resolution of y_i obtained by using Eq. (31), bases ϕ_v , ϕ_m , and ϕ_σ are determined in this section. To simplify the algorithm, let N_v , N_m , and N_σ be set to N and let N be given in advance. Let N_{v_i} , N_{m_i} , and N_{σ_i} be set to N_i , let ϕ_v , ϕ_m , and ϕ_σ be set to $\phi \stackrel{\text{def}}{=} {}^t(\phi_0, \dots, \phi_N)$, and let $\phi_i \stackrel{\text{def}}{=} K(\mathbf{x}, \mathbf{k})$ as shown in Eq. (A·2) in Appendix A. Here, $K(\mathbf{x}, \mathbf{k})$ is a multi-dimensional orthonormal basis defined by Eq. (A·1), \mathbf{k} is the index vector, and the relationship between i , the index of the basis, and \mathbf{k} is obtained by Eq. (A·3).

Once N and $K(\mathbf{x}, \mathbf{k})$ are given, ϕ is determined by setting N_i . As mentioned in Sect. 3.5, to obtain a large discount return, it is reasonable to assign a larger value to N_i if $res_{y_{u_i}}$ is large. Furthermore, Eq. (A·4), which shows the relationship between N and N_i , should hold. Thus, to satisfy these requirements, we determine N_i by the following equation:

$$N_i = \text{int}(\alpha res_{y_{u_i}}), \quad (33)$$

where $\text{int}(\cdot)$ denotes the integer part, and α is set so that $|N - \prod_{i=1}^{d_y} (N_i + 1) - 1|$ is minimum for the given N and $res_{y_{u_i}}$.

Let ϕ be the basis when the degree of expansion is N_j , and ϕ' be the basis when the degree of expansion is N'_j . When the degree of expansion changes from N_j to N'_j , the basis also changes from ϕ to ϕ' . Thus, new coefficients c' , $\xi'_{m_{i'}}$, and $\xi'_{\sigma_{i'}}$ for the new bases have to take the values of the old ones, c , ξ_{m_i} , and ξ_{σ_i} , for learning to proceed smoothly. That is, $c'_{i'}$ is set to c_i , $\xi'_{m_{i'}}$ is set to ξ_{m_i} , and $\xi'_{\sigma_{i'}}$ is set to ξ_{σ_i} , if $\phi'_{i'}$ is equal to ϕ_i . Further details are described in Appendix A.

4. Expansion to Multiple Agents

Let us consider an environment that L agents control, and in which the agents cannot always observe all the state variables and the control inputs. In this section, the agent described by Algorithm 2 is extended to multiple agents. Let \hat{x}_i be the i th state variable of the environment, $\hat{\mathbf{x}} \stackrel{\text{def}}{=} {}^t(\hat{x}_1, \dots, \hat{x}_{\hat{d}_x}) \in \hat{\mathbf{D}}_x$, $\hat{\mathbf{D}}_x \stackrel{\text{def}}{=} \{\hat{x}_i | \hat{x}_{\min_i} \leq \hat{x}_i \leq \hat{x}_{\min_i} + \hat{D}_{x_i}, 1 \leq i \leq \hat{d}_x\}$, \hat{v}_i be the i th control input of the environment, $\hat{\mathbf{v}} \stackrel{\text{def}}{=} {}^t(\hat{v}_1, \dots, \hat{v}_{\hat{d}_v})$, $\hat{\mathbf{D}}_v \stackrel{\text{def}}{=} \{\hat{v}_i | \hat{v}_{\min_i} \leq \hat{v}_i \leq \hat{v}_{\min_i} + \hat{D}_{v_i}, 1 \leq i \leq \hat{d}_v\}$, $r^{(\ell)}$ be the immediate reward that the ℓ th agent receives, $\mathbf{x}_t^{(\ell)}$ be the subset of $\hat{\mathbf{x}}_t$ that the ℓ th agent can observe, $\mathbf{v}_t^{(\ell)}$ be the subset of $\hat{\mathbf{v}}_t$ that the ℓ th agent can observe, and $\mathbf{u}_t^{(\ell)}$ be the action determined by the ℓ th agent. To express the relationship between $\mathbf{x}_t^{(\ell)}$ and $\hat{\mathbf{x}}_t$, and that between $\mathbf{v}_t^{(\ell)}$ and $\hat{\mathbf{v}}_t$, let us introduce state observation matrix $G_s^{(\ell)}$ and action observation matrix $G_v^{(\ell)}$. By using $G_s^{(\ell)}$ and $G_v^{(\ell)}$, the ℓ th agent observes $\mathbf{x}_t^{(\ell)}$ and

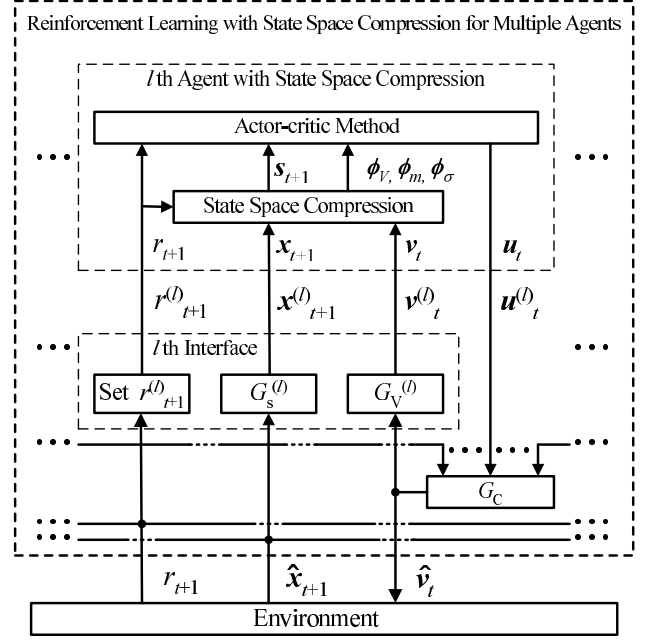


Fig. 2 Structure of reinforcement learning with state space compression for multiple agents.

$\mathbf{v}_t^{(\ell)}$ as follows:

$$\begin{aligned} \mathbf{x}_t^{(\ell)} &= G_s^{(\ell)} \hat{\mathbf{x}}_t, \\ \mathbf{v}_t^{(\ell)} &= G_v^{(\ell)} \hat{\mathbf{v}}_t. \end{aligned} \quad (34)$$

Actions $\mathbf{u}_t^{(1)}, \dots, \mathbf{u}_t^{(L)}$ are combined into $\hat{\mathbf{v}}_t$ by

$$\hat{\mathbf{v}}_t = G_c {}^t(\mathbf{u}_t^{(1)}, \dots, \mathbf{u}_t^{(L)}) + \mathbf{g}_c. \quad (35)$$

Here, G_c and \mathbf{g}_c are determined from v_{\min_i} , \hat{D}_{v_i} , $u_{\min_i}^{(\ell)}$, and $D_{u_i}^{(\ell)}$. Immediate reward $r_t^{(\ell)}$ for the ℓ th agent is set based on r_t .

The structure and algorithm of reinforcement learning with state space compression for multiple agents are shown in Fig. 2 and given by Algorithm 4, respectively.

Algorithm 4: Reinforcement learning with state space compression for multiple-agents.

- (4-1) Initialization-1: Set $\mathbf{c}_0^{(\ell)}$, $\xi_{m0}^{(\ell)}$, $\xi_{\sigma0}^{(\ell)}$, $n_{\text{MA}}^{(\ell)} = 0$, and $\text{trg}_{\text{MA}}^{(\ell)} = 0$ for $1 \leq \ell \leq L$.
- (4-2) Initialization-2: Set $t = 0$, $\hat{\mathbf{x}}_0$, and $\hat{\mathbf{v}}_0$.
- (4-3) Environment: Receive $\hat{\mathbf{v}}_t$, and output $\hat{\mathbf{x}}_{t+1}$ to the interface.
- (4-4) Interface: Set $\mathbf{x}_{t+1}^{(\ell)}$, $\mathbf{x}_t^{(\ell)}$, and $\mathbf{v}_t^{(\ell)}$ by using Eq. (34) for $1 \leq \ell \leq L$, and set $r_{t+1}^{(\ell)}$ for $1 \leq \ell \leq L$.
- (4-5) Agent with State Space Compression: Compute $\mathbf{u}_{t+1}^{(\ell)}$ by using Step (2-4) in Algorithm 2 for $1 \leq \ell \leq L$.
- (4-6) G_c : Compute $\hat{\mathbf{v}}_{t+1}$ by using Eq. (35).

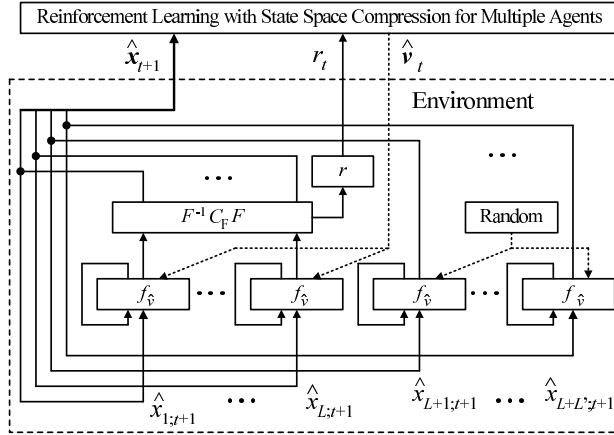


Fig. 3 Simple environment: Linear combination of logistic maps.

(4-7) Set $t = t + 1$. If an episode finished, go to Step (4-2), else go to Step (4-3).

5. Performance Evaluation

Consider the simple example environment shown in Fig. 3, which is a linear combination of $L + L'$ logistic maps [8]. Here, $d_u^{(\ell)} = 1$ for $1 \leq \ell \leq L$, $\hat{d}_v = d_v = L$, $\hat{d}_x = d_x = L + L'$, $0 \leq \hat{x}_{\ell;t} \leq 1$ for $1 \leq \ell \leq L + L'$, I_ℓ is the $(\ell \times \ell)$ -identity matrix, $G_s^{(\ell)} = I_{L+L'}$, $G_v^{(\ell)} = I_L$ for $1 \leq \ell \leq L$, i.e., $\mathbf{x}_t^{(\ell)} = \hat{\mathbf{x}}_t$ and $\mathbf{v}_t^{(\ell)} = \hat{\mathbf{v}}_t$ for $1 \leq \ell \leq L$, $G_c = I_L$ and $\mathbf{g}_c = 0$, i.e., $\hat{\mathbf{v}}_t = {}^t(u_1^{(1)} \ t, \dots, u_1^{(L)} \ t)$, and the logistic map is defined by the following equation [8]:

$$\begin{aligned} x' &= f_v(x) \\ &\stackrel{\text{def}}{=} vx(1-x). \end{aligned} \quad (36)$$

States $\hat{x}'_{1;t}, \dots, \hat{x}'_{L;t}$ interact with each other, and $\hat{\mathbf{x}}_{t+1}$ is obtained by the following equation:

$$\hat{x}'_{\ell;t} = \overbrace{f_{\hat{v}_{\ell;t}} \cdots f_{\hat{v}_{\ell;t}}}^{n_{\text{map}} \text{ times}}(\hat{x}_{\ell;t}) \quad \text{for } 1 \leq \ell \leq L + L', \quad (37)$$

$$\begin{cases} {}^t(\hat{x}'_{1;t+1}, \dots, \hat{x}'_{L;t+1}) \\ \quad = F^{-1} C_F F^t(\hat{x}'_{1;t}, \dots, \hat{x}'_{L;t}), \\ {}^t(\hat{x}'_{L+1;t+1}, \dots, \hat{x}'_{L+L';t+1}) \\ \quad = {}^t(\hat{x}'_{L+1;t}, \dots, \hat{x}'_{L+L';t}), \end{cases}$$

where $n_{\text{map}} = 3$, $3.4 \leq \hat{v}_{\ell;t} \leq 4$ for $1 \leq \ell \leq L + L'$, F is a matrix that performs discrete Fourier transform, and $C_F \stackrel{\text{def}}{=} \text{diag}[1.0, 0.9, 0.81, 0.73, 0.66, 0.59, 0.53, 0.48, \dots]$, which denotes a low pass filter that reduces the power spectrum of $\hat{x}'_{1;t}, \dots, \hat{x}'_{L;t}$ as the frequency increases.

The ℓ th agent determines action $u_1^{(\ell)} = \hat{v}_{\ell;t}$ by using the reinforcement learning described in Algorithm 4, and $\hat{v}_{L+1;t}, \dots, \hat{v}_{L+L';t}$ are set so that they independently obey a uniform distribution. Let $\text{spc}_{0;t}, \dots, \text{spc}_{L-1;t}$ be the power spectrum of

$\hat{x}'_{1;t}, \dots, \hat{x}'_{L;t}$, and immediate reward $r_{t+1}^{(\ell)}$ for $1 \leq \ell \leq L$ be equal to r_{t+1} defined by

$$r_{t+1} \stackrel{\text{def}}{=} (\Delta \text{spc}_{k;t+1})^2 - \sum_{j=1, j \neq k}^L (\Delta \text{spc}_{j;t+1})^2, \quad (38)$$

where the power spectrum denotes the square of the Fourier coefficients obtained by $F^t(\hat{x}'_{1;t}, \dots, \hat{x}'_{L;t})$, $\Delta \text{spc}_{j;t+1} \stackrel{\text{def}}{=} \text{spc}_{j;t+1} - \text{spc}_{j;t}$, k is set to 2. Thus, r_{t+1} increases when $\hat{x}'_{1;t}, \dots, \hat{x}'_{L;t}$ move so that the motion in $\text{spc}_{2;t}$ is large. Because $\hat{x}'_{L+1;t}, \dots, \hat{x}'_{L+L';t}$ are not related to the immediate reward, they are disturbance states that disturb the 1st to the L th agent.

The effectiveness of the state space compression algorithm described in Sect. 3 and the ability of Algorithm 4 to construct a feature space without being disturbed by the disturbance states were demonstrated by evaluating Algorithm 4 for various values for the degree of expansion, N , and the number of disturbance states, L' , and for several combinations of $\boldsymbol{\gamma} \stackrel{\text{def}}{=} (\gamma_1, \gamma_2, \gamma_3, \gamma_4)$ and $\tilde{\boldsymbol{\gamma}} \stackrel{\text{def}}{=} (\tilde{\gamma}_1, \tilde{\gamma}_2, \tilde{\gamma}_3, \tilde{\gamma}_4)$. Note that the parameters for each agent in Algorithm 4 have the same values, and superscript (ℓ) for variables and parameters is omitted in the following to simplify the explanation. The combinations of $\boldsymbol{\gamma}$ and $\tilde{\boldsymbol{\gamma}}$ evaluated were determined as follows. (1) $\tilde{\gamma}_1$: It is reasonable to construct a feature space based on the range of the i th principal component, y_i . In other words, degree of expansion N_i for y_i should be proportional to $\sigma(y_i)$. Thus, $\tilde{\gamma}_1$ is set to 1. (2) γ_2 : Because it is clear that the difference in D_{v_j} should be normalized, γ_2 is set to 1. Note that γ_2 does not affect the performance shown below, because $d_u = 1$ in the environment described in Eq. (37). (3) γ_3 and γ_4 : It is important to distinguish the states that contribute to the reward from the states that disturb the control algorithm. To do this, the effect of $v_{j;t}$ on $x_{i;t+1}$ and those of $x_{i;t}$ and $x_{i;t+1}$ on r_{t+1} are evaluated. The former is expressed by $|\beta_{v_{ij}}|^{\gamma_3}$, and the latter is expressed by $(|b_{xi}| + |b_{x'i}|)^{\gamma_4}$. Thus, Algorithm 4 is evaluated for $\gamma_3 = 0$ or 1, and $\gamma_4 = 0$ or 1. (4) γ_1 , $\tilde{\gamma}_2$, $\tilde{\gamma}_3$, and $\tilde{\gamma}_4$: Because $\sigma(x_i)$ and $\sigma(y_i)$ express the same property (i.e., the standard deviation in the state) of the environment, γ_1 is set to 0 so that an agent does not take the standard deviation into account more than once. For the same reason, $\tilde{\gamma}_2$, $\tilde{\gamma}_3$, and $\tilde{\gamma}_4$ are set to 0. Therefore, $\tilde{\boldsymbol{\gamma}}$ is set to $(1, 0, 0, 0)$, and $\boldsymbol{\gamma}$ is set to $(0, 0, 0, 0)$, $(0, 1, 1, 0)$, $(0, 1, 0, 1)$, or $(0, 1, 1, 1)$, where the motions in Algorithm 4 for $(0, 0, 0, 0)$ and $(0, 1, 0, 0)$ are the same in the environment of Eq. (37), and $(0, 0, 0, 0)$ means that $\text{res}_{x_{u1}} = \dots = \text{res}_{x_{uL+L'}} = 1$. The other parameters for reinforcement learning are set as follows: $T_S = 5 \times 10^6$, $T_{RB} = 10^3$, $\nu_{Tr} = 0.3$, $\nu_{DR} = 0.9$, and $\zeta_v = \zeta_m = \zeta_\sigma = 10^{-4}$.

The effect of N and L on Algorithm 4 was estimated by evaluating the change in immediate reward r_t when $\boldsymbol{\gamma} = (0, 0, 0, 0)$ and $L' = 0$. The performance of

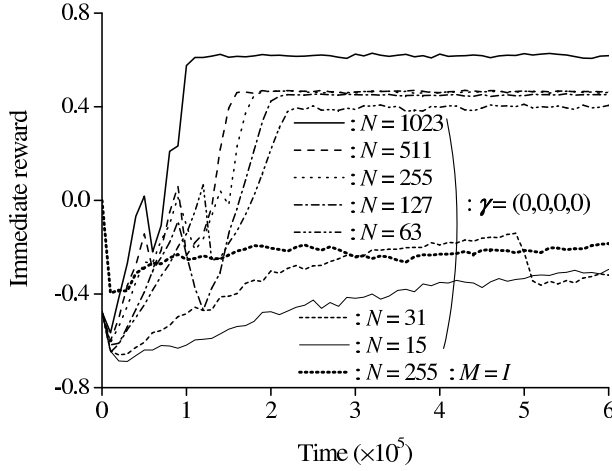


Fig. 4 Change in immediate reward for various values of degree of expansion after starting reinforcement learning ($L = 8$).

the state space compression in Algorithm 4 was demonstrated by also evaluating an algorithm that does not use the state space compression, in which $M = I$ and N_i was given in advance ($N_1 = \dots = N_L = 1$; i.e., $N = 255$ when $L = 8$). Figure 4 shows the change in r_t for $L = 8$ after starting reinforcement learning; that is, the time in Figure 4 indicates $t - T_S$. The figure shows that Algorithm 4 with $\gamma = (0, 0, 0, 0)$ provides a sufficiently large r_t when N is at least 63. On the other hand, Algorithm 4 without state space compression (i.e., $M = I$) provides a lower value of r_t although $N (= 255)$ is larger than 63. Therefore, the state space compression method described in Sect. 3 is effective. Figure 5 shows the effect of N and L on r_t , where r_t is normalized by the maximum value of r_t for each value of L . The changes in the normalized reward for $L = 4, 8$, and 16 are almost the same as shown in this figure. This means that the performance is not always affected by L , the dimension of the system. The factors affecting the performance could be the complexity of the environment defined by n_{map} and C_F .

The effect of the number of disturbance states on immediate reward r_t is shown in Fig. 6, where $N = 255$ and $L = 8$. This figure shows that the effect of γ_3 and γ_4 becomes noticeable as L' , the number of disturbance states, increases, and that the largest r_t is obtained when $\gamma_3 = \gamma_4 = 1$ for each value of L' . Thus, $|\beta_{v_{ij}}|$ and $|b_{x_i}| + |b_{x'_i}|$ are effective to remove the effect of the disturbance states, and we can conclude that Algorithm 4 with $\gamma = (0, 1, 1, 1)$ constructs a feature space effectively and controls the environment without suffering the disturbance states.

The motion in state $\hat{x}_{\ell;t}$ was demonstrated by measuring the power spectrum and the change in state $\hat{x}_{1;t}, \dots, \hat{x}_{L;t}$ for $M = I$ and $\gamma = (0, 1, 1, 1)$ after r_t was stable and saturated. The results are shown in Figs. 7 to 9 for $N = 255$, $L = 8$, and $L' = 0$. The origin of the time in Figs. 8 and 9 indicates the time

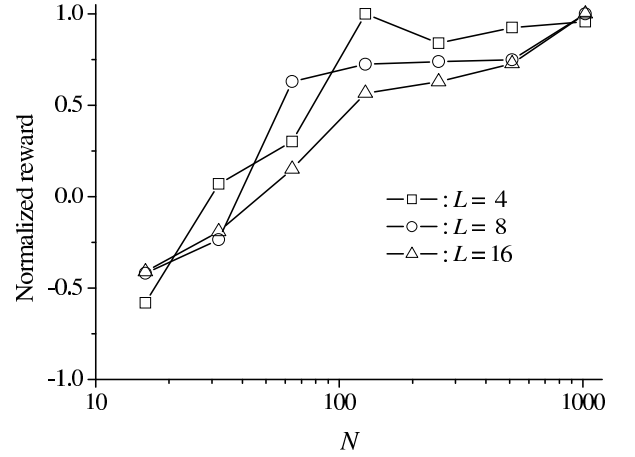


Fig. 5 The effect of N and L on normalized reward.

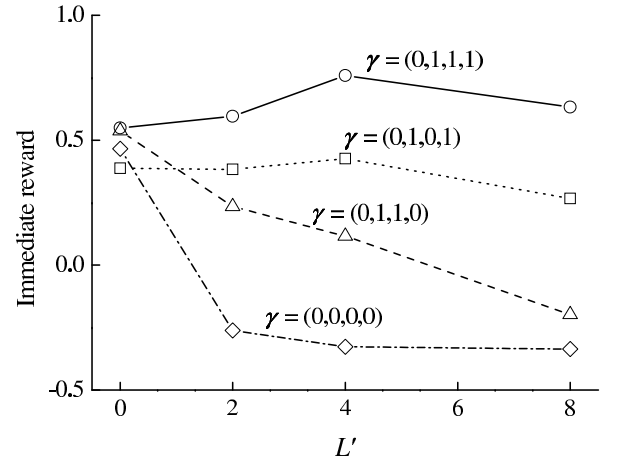


Fig. 6 Effect of number of disturbance states on immediate reward ($N = 255$ and $L = 8$).

when measurement started. These figures show that spc_2 in Eq. (38) is extracted and Algorithm 4 works well when $\gamma = (0, 1, 1, 1)$, although $\hat{x}_{1;t}, \dots, \hat{x}_{L;t}$ moves almost randomly when $M = I$.

All combinations of γ and $\tilde{\gamma}$ were not evaluated, and it was assumed that all states and control inputs in the environment could be observed. To better clarify the effectiveness of Algorithm 4, a more extensive investigation will be conducted.

6. Conclusion

A state space compression method based on multivariate analysis was developed and applied to reinforcement learning for high-dimensional continuous state spaces. It can autonomously construct a feature space without preliminary knowledge of the environment. An example synchronization problem for multiple logistic maps was solved, demonstrating that the state space compression method exhibits high performance without suffering from disturbance states. Because the basis of the

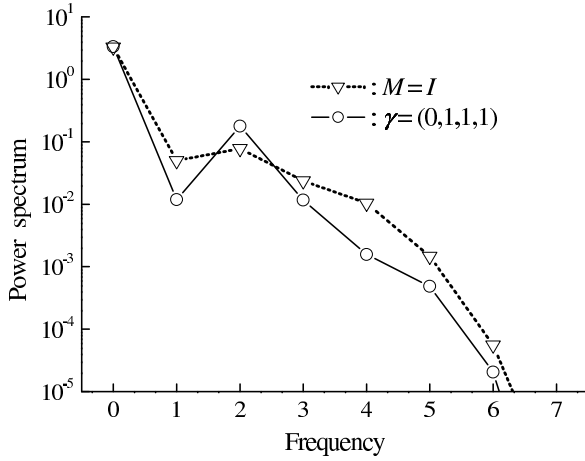


Fig. 7 Power spectrum of $\hat{x}_1, \dots, \hat{x}_8$ ($N = 255$, $L = 8$, and $L' = 0$).

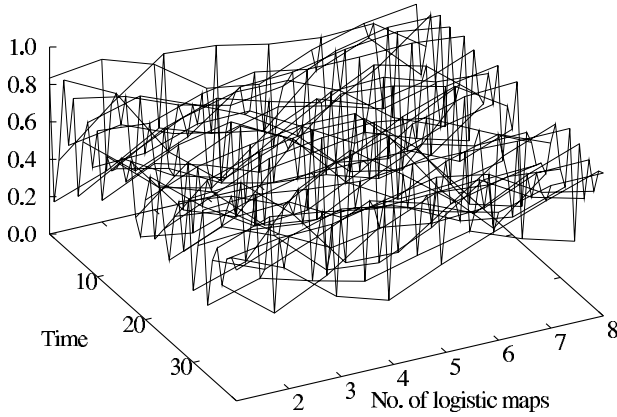


Fig. 8 Motion in $\hat{x}_1, \dots, \hat{x}_8$ for $M = I$ ($N = 255$, $L = 8$, and $L' = 0$).

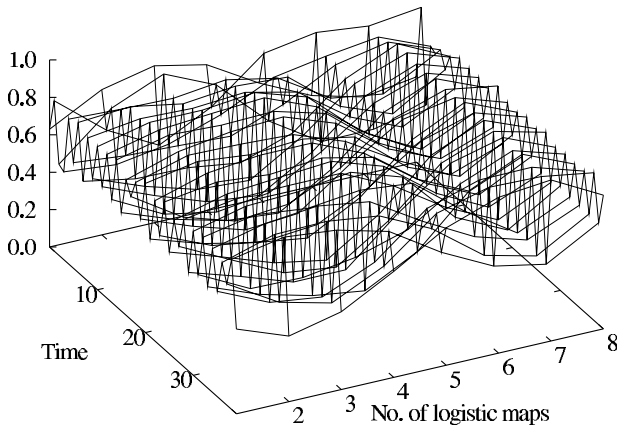


Fig. 9 Motion in $\hat{x}_1, \dots, \hat{x}_8$ for $\gamma = (0, 1, 1, 1)$ ($N = 255$, $L = 8$, and $L' = 0$).

feature space is orthogonal, the mathematical structure of the system in the feature space is clear, and it is possible to mathematically analyze the dynamics and stability in a moment vector space. The mathematical

analysis of the dynamics will be reported in a future.

Appendix A: Orthonormal Basis for Function Approximation

Let $\mathbf{x} \stackrel{\text{def}}{=} {}^t(x_1, \dots, x_{d_x})$ and $\mathcal{D}_x \stackrel{\text{def}}{=} \{x_j | x_{\min_j} \leq x_j \leq x_{\min_j} + D_{x_j}, 1 \leq j \leq d_x\}$. In this appendix, the orthonormal bases [12] for $\mathbf{x} \in \mathcal{D}_x$ are summarized.

Let $\mathbf{k} \stackrel{\text{def}}{=} {}^t(k_1, \dots, k_{d_x})$ be the index vector of the Fourier coefficient and $h(\mathbf{k})$ be the Fourier coefficient for index vector \mathbf{k} . The Fourier series expansion for function $f(\mathbf{x})$ is defined by the following equation [12]:

$$f(\mathbf{x}) = \sum_{\mathbf{k} \in \mathcal{Z}} h(\mathbf{k}) K(\mathbf{x}, \mathbf{k}),$$

$$h(\mathbf{k}) \stackrel{\text{def}}{=} \int_{\mathcal{D}_x} f(\mathbf{x}) K^*(\mathbf{x}, \mathbf{k}) d\mathbf{x},$$

where $\mathcal{Z} \stackrel{\text{def}}{=} \{k_j | 0 \leq k_j \leq N_j, 1 \leq j \leq d_x\}$, N_j is the degree of expansion of x_j , $h(\mathbf{k})$ is the Fourier coefficient, superscript $*$ denotes a complex conjugate, and $\{K(\mathbf{x}, \mathbf{k})\}$ is a multi-dimensional orthonormal basis defined by the following equation:

$$K(\mathbf{x}, \mathbf{k}) \stackrel{\text{def}}{=} \prod_{j=1}^{d_x} K_j(x_j, k_j). \quad (\text{A.1})$$

Here, $\{K_j(x_j, k_j)\}$ is one-dimensional orthonormal basis, $K_j(x_j, k_j)$ is set to a trigonometrical function defined by

$$K_j(x_j, k_j) \stackrel{\text{def}}{=} \begin{cases} \sqrt{\frac{1}{D_{x_j}}} & \text{for } k_j = 0 \\ \sqrt{\frac{2}{D_{x_j}}} \sin\left(\frac{k_j + 1}{2} \omega_{0_j}(x_j - x_{\min_j})\right) & \text{for } k_j \in \mathcal{N}_o, \\ \sqrt{\frac{2}{D_{x_j}}} \cos\left(\frac{k_j}{2} \omega_{0_j}(x_j - x_{\min_j})\right) & \text{for } k_j \in \mathcal{N}_e \end{cases}$$

$\omega_{0_j} \stackrel{\text{def}}{=}} 2\pi/D_{x_j}$, \mathcal{N}_e denotes a set of even numbers, and \mathcal{N}_o denotes a set of odd numbers.

Let us define basis $\phi(\mathbf{x}) = {}^t(\phi_0(\mathbf{x}), \dots, \phi_N(\mathbf{x}))$ as

$$\phi_i(\mathbf{x}) \stackrel{\text{def}}{=} K(\mathbf{x}, \mathbf{k}), \quad (\text{A.2})$$

where i is the index of the basis, N is the degree of expansion of \mathbf{x} , and $N+1$ is the dimension of the feature space with the basis. When $\mathcal{Z} \stackrel{\text{def}}{=} \{k_j | 0 \leq k_j \leq N_j, 1 \leq j \leq d_x\}$, that is, \mathcal{Z} is a cube, the index of the basis is obtained from index vector \mathbf{k} as follows:

$$i = \sum_{j=1}^{d_x} k_j \prod_{\nu=j+1}^{d_x} N_\nu. \quad (\text{A.3})$$

From the degree of expansion of x_j , we can obtain the degree of expansion of \mathbf{x} as follows:

$$N = \prod_{j=1}^{d_x} (N_j + 1) - 1. \quad (\text{A}\cdot 4)$$

Given that the degree of expansion of x_j is N_j , suppose that the basis, the index of the basis, and the domain of \mathbf{k} are given by ϕ , $i(\mathbf{k})$, and \mathcal{Z} , respectively. When the degree of expansion of x_j is N'_j , ϕ' , $i'(\mathbf{k})$, and \mathcal{Z}' are also given. As mentioned in Sect. 3.7, when the degree of expansion changes from N_j to N'_j , the new coefficients, $c'_{i'(\mathbf{k})}$, $\xi'_{m_{i'(\mathbf{k})}}$, and $\xi'_{\sigma_{i'(\mathbf{k})}}$, for the new bases take the values of the old ones, $c_{i(\mathbf{k})}$, $\xi_{m_{i(\mathbf{k})}}$, and $\xi_{\sigma_{i(\mathbf{k})}}$, as follows:

$$\begin{aligned} c'_{i'(\mathbf{k})} &= c_{i(\mathbf{k})} & \text{for } \forall \mathbf{k} \in \mathcal{Z}', \\ \xi'_{m_{i'(\mathbf{k})}} &= \xi_{m_{i(\mathbf{k})}} & \text{for } \forall \mathbf{k} \in \mathcal{Z}', \\ \xi'_{\sigma_{i'(\mathbf{k})}} &= \xi_{\sigma_{i(\mathbf{k})}} & \text{for } \forall \mathbf{k} \in \mathcal{Z}'. \end{aligned} \quad (\text{A}\cdot 5)$$

If $\mathbf{k} \in \mathcal{Z}'$ and $\mathbf{k} \notin \mathcal{Z}$, that is, there is no $\phi_{i(\mathbf{k})}$ corresponding to $\phi'_{i'(\mathbf{k})}$, $c'_{i'(\mathbf{k})}$, $\xi'_{m_{i'(\mathbf{k})}}$, and $\xi'_{\sigma_{i'(\mathbf{k})}}$ are set to 0.

Appendix B: TD Learning Based on Steepest Descent Method

TD learning [1] based on the steepest descent method is summarized in this appendix. Let $\mathbf{x} \stackrel{\text{def}}{=} {}^t(x_1, \dots, x_{d_x})$, $\mathcal{D}_{\mathbf{x}} \stackrel{\text{def}}{=} \{x_i | x_{\min_i} \leq x_i \leq x_{\min_i} + D_{x_i}, 1 \leq i \leq d_x\}$, $p(\mathbf{x})$ be a PDF of \mathbf{x} , $g(\mathbf{x}, \mathbf{c})$ be an approximation of function $g(\mathbf{x})$, and \mathbf{c} be the parameter used by the approximation. Parameter \mathbf{c} is determined so as to minimize $MSE(\mathbf{c})$ defined by

$$MSE(\mathbf{c}) \stackrel{\text{def}}{=} \int_{\mathbf{x} \in \mathcal{D}_{\mathbf{x}}} p(\mathbf{x})(g(\mathbf{x}) - g(\mathbf{x}, \mathbf{c}))^2 d\mathbf{x}. \quad (\text{A}\cdot 6)$$

Let us now derive parameter \mathbf{c} by using the following steepest descent method [13].

Algorithm 5: Steepest Descent Method.

- (1) Set initial value \mathbf{c}_0 and step size $\zeta > 0$.
- (2) Set number of iterations $t = 0$.
- (3) $\mathbf{c}_{t+1} = \mathbf{c}_t - \zeta \nabla_{\mathbf{c}} MSE(\mathbf{c}_t)$.
- (4) $t = t + 1$.
- (5) Go to Step (3).

Here, \mathbf{c}_t is the value of vector \mathbf{c} at t .

To realize Step (3) in Algorithm 5, let us make the following assumption.

Assumption 1: PDF $p(\mathbf{x})$ is the same as the PDF when Algorithm 5 has been performed.

By applying Assumption 1 and Eq. (A·6) to Step (3) in Algorithm 5, and by replacing ζ with $\zeta/2$, Step (3) in Algorithm 5 can be described as follows:

$$\begin{aligned} \mathbf{c}_{t+1} &= \mathbf{c}_t - (\zeta/2) \nabla_{\mathbf{c}} MSE(\mathbf{c}_t) \\ &= \mathbf{c}_t - (\zeta/2) \nabla_{\mathbf{c}} (g(\mathbf{x}_t) - g(\mathbf{x}_t, \mathbf{c}_t))^2 \\ &= \mathbf{c}_t + \zeta (g(\mathbf{x}_t) - g(\mathbf{x}_t, \mathbf{c}_t)) \nabla_{\mathbf{c}} g(\mathbf{x}_t, \mathbf{c}_t). \end{aligned} \quad (\text{A}\cdot 7)$$

Let $\hat{g}(\mathbf{x})$ be the estimated value of $g(\mathbf{x})$. When $g(\mathbf{x}_t)$ is unknown but $\hat{g}(\mathbf{x})$ can be obtained, Step (3) of Algorithm 5 is obtained by replacing $g(\mathbf{x})$ with $\hat{g}(\mathbf{x})$ in Eq. (A·7) as follows:

$$\mathbf{c}_{t+1} = \mathbf{c}_t + \zeta (\hat{g}(\mathbf{x}_t) - g(\mathbf{x}_t, \mathbf{c}_t)) \nabla_{\mathbf{c}} g(\mathbf{x}_t, \mathbf{c}_t). \quad (\text{A}\cdot 8)$$

Consider that $g(\mathbf{x}_t)$ is unknown and that we cannot obtain $\hat{g}(\mathbf{x})$; that is, we can obtain only an estimated value, $\tilde{g}(\mathbf{x})$, of $g(\mathbf{x})$ with unknown accuracy. In this case, introducing TD error $\varepsilon_{\text{TD}t}$ into Eq. (A·7) and replacing $g(\mathbf{x})$ with $\tilde{g}(\mathbf{x})$ enable Step (3) of Algorithm 5 to be expressed by the following equation:

$$\mathbf{c}_{t+1} = \mathbf{c}_t + \zeta \varepsilon_{\text{TD}t} (\tilde{g}(\mathbf{x}_t) - g(\mathbf{x}_t, \mathbf{c}_t)) \nabla_{\mathbf{c}} g(\mathbf{x}_t, \mathbf{c}_t). \quad (\text{A}\cdot 9)$$

References

- [1] R. S. Sutton and A. G. Barto, Reinforcement Learning, MIT Press, USA, 1998.
- [2] Baird, L. C. and Klopff, A. H., "Reinforcement learning with high-dimensional, continuous actions," Technical Report WL-TR-93-1147, 1993.
- [3] W. T. B. Uther and M. M. Veloso, "Tree based discretization for continuous state space reinforcement learning," Proceedings of AAAI-98, pp. 769–774, July, 1998.
- [4] A. J. Smith, "The applications of the self-organising map to reinforcement learning," Neural Networks archive, vol. 15, Issue 8-9, pp. 1107–1124, Oct. 2002.
- [5] T. Dietterich, "Hierarchical reinforcement learning with the MAXQ value function decomposition," Journal of Artificial Intelligence Research, vol. 13, pp. 227–303, 2000.
- [6] A. G. Barto and S. Mahadevan, "Recent advances in hierarchical reinforcement learning," Discrete event dynamic systems, vol. 13, Issue 4, pp. 341–379, Oct., 2003.
- [7] H. Satoh, "Approximation and Analysis of Non-linear Equations Based on Moment Vector Equation," Tech. Rept. IEICE. NLP2005-1, pp. 1-6, Jan. 2005.
- [8] E. Ott, Chaos in Dynamical Systems Second edition, Cambridge, UK, 2002.
- [9] I. H. Witten, "An adaptive optimal controller for discrete-time Markov environments," Information and Control, vol. 34, pp. 286–295, 1977.
- [10] A. G. Barto, R. S. Sutton, and C. W. Anderson, "Neuron-like elements that can solve difficult learning control problems," IEEE Trans. on Systems Man and Cybernetics, vol. 13, pp. 835–846, 1983.
- [11] R. A. Johnson, D. W. Wichern, Applied Multivariate Statistical Analysis, 5th ed, Pearson Education (Prentice Hall USA), 2001.
- [12] I. N. Bronshtein and K. A. Semendyayev, Handbook of Mathematics, Springer-Verlag, Berlin, 1997.
- [13] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, Numerical Recipes in C 2nd Edition, Cambridge University Press, UK, 1992.