

# Vision-based Learning and Development for Emergence of Robot Behaviors

Minoru ASADA, Koh HOSODA, and Sho'ji SUZUKI

Dept. of Adaptive Machine Systems, Graduate School of Engineering

Osaka University, Suita, Osaka 565 (Japan)

email: asada@ams.eng.osaka-u.ac.jp

**Abstract:** This paper focuses on two issues on learning and development; a problem of state-action space construction, and a scaling-up problem. The former is mainly related to sensory-motor mapping and its abstraction, and we show two our methods for the state and action space construction for reinforcement learning. For the latter issue, we attempt to define the environmental complexity based on the relationships between observations and self motions. Based on this view, we introduce a method which can cope with the complexity of multi-agent environment by a combination of a state vector estimation process and a reinforcement learning process based on the estimated vectors. As example tasks in our work, we adopt the domain of soccer robots, RoboCup [1]. Computer simulations and real robot experiments are given.

## 1. Introduction

The ultimate goal of our research is to design the fundamental internal structure inside physical entities having their bodies (robots) which can emerge complex behaviors through the interactions with their environments. In order to emerge the intelligent behaviors, physical bodies have an important role of bringing the system into *meaningful* interaction with the physical environment – complex, uncertain, but with automatically consistent set of natural constraints. This facilitates the correct agent design, learning from the environment, and rich meaningful agent interaction. The meanings of “having a physical body” can be summarized as follows:

1. Sensing and acting capabilities are not separable, but tightly coupled.
2. In order to accomplish the given tasks, the sensor and actuator spaces should be abstracted under the resource bounded conditions (memory, processing power, controller etc.).
3. The abstraction depends on both the fundamental embodiments inside the agents and the experiences (interactions with their environments).
4. The consequences of the abstraction are the agent-based subjective representation of the environment, and its evaluation can be done

by the consequences of behaviors.

5. In real world, both inter-agent and agent-environment interactions are asynchronous, parallel and arbitrarily complex. The agent should cope with increasing complexity of the environment to accomplish the given task at hand.

In this paper, we focus on two issues on learning and development; a problem of state-action space construction, and a scaling-up problem. The former is mainly related to 2 and 3, and we show two our methods for the state and action space construction for reinforcement learning. One is based on an off-line learning method [2] and the other on-line one [3].

The latter issue is closely related to 4 and 5, and we attempt to define the environmental complexity based on the relationships between observations and self motions. Based on this view, we introduce a method which can cope with the complexity of multi-agent environment by a combination of a state vector estimation process and a reinforcement learning process based on the estimated vectors [4].

As example tasks in our work, we adopt the domain of soccer robots, RoboCup, which is an attempt to foster intelligent robotics research by providing a standard problem where a wide range of technologies can be integrated and examined [1].

The remainder of this article is structured as follows. We first give an explanation of the problem of state-action space construction along with our real robot experiments in the context of reinforcement learning. Then, we show our method to cope with more complicated tasks in multi-agent environment. Finally, we give a conclusion.

## 2. A Problem of State-Action Space Construction

Reinforcement learning [5, 6] has been receiving increased attention as a method for robot learning with little or no *a priori* knowledge and higher capability of reactive and adaptive behaviors. In such robot learning methods, a robot and an environment are generally modeled by two synchronized finite state automata interacting in a discrete time cyclical processes. The robot senses the current state of the environment and selects an action. Based on the state and the action, the environment makes a transition to a new state and generates a reward that is passed back to the robot. Through these interactions, the robot learns a purposive behavior to achieve a given goal.

To apply robot learning methods such as reinforcement learning to real robot tasks, we need a well-defined state-action space by which the robot learns to select an adequate action for the current state to accomplish the task at hand. Traditional notions of state and action in the existing applications of the reinforcement learning schemes fit nicely into deterministic state transition models (e.g. one action is forward, backward, left, or right, and the states are encoded by the locations of the agent). However, it seems difficult to apply such deterministic state transition models to real robot tasks. In real world, everything changes asynchronously [7]. Therefore, the construction of state-action space is one of the most important issues in robot learning.

Generally, the design of the state-action space in which the necessary and sufficient information to accomplish a given task should be included depends on the capabilities of agent sensing and acting. The abstraction process from sensory information to a state seems to depend on the process from motor commands to an action, and *vice versa*. This resembles the well-known “chicken and egg problem” that is difficult to be solved (see Figure 1).

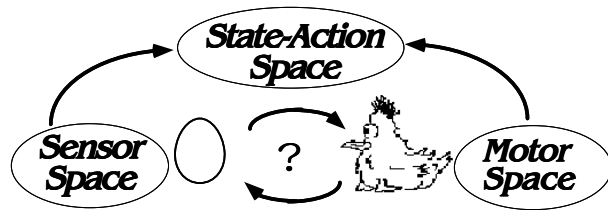


Fig. 1. The inter-dependence between sensor and motor spaces from a viewpoint of state-action space construction

### 2.1 An Off-Line Learning Method

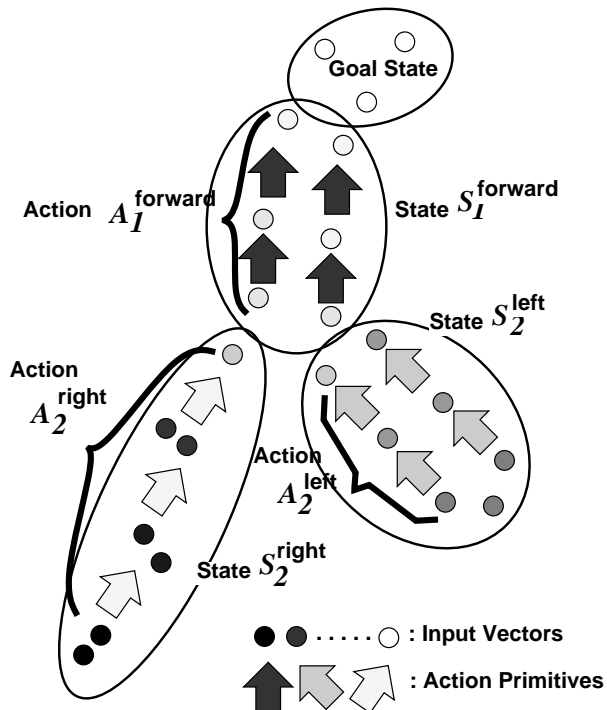
Basic ideas of our first approach to cope with this problem are:

- we define an action primitive as a motor command to be executed during a fixed time interval, and an input vector as sensory data of the consequence of the action primitive, and
- we define a state as a cluster of input vectors from which the robot can reach the goal state (or the state already obtained) by a sequence of one kind action primitive regardless of its length, and one action as this sequence of action primitives.

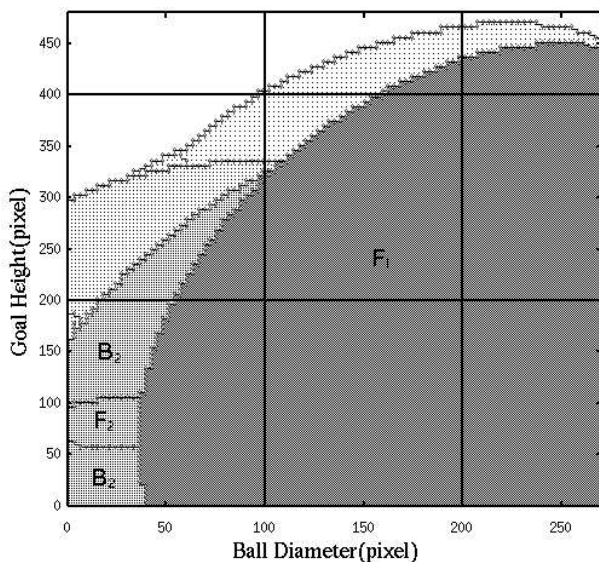
Figure 2 (a) shows the basic idea of the state-action space construction. The initial state space consisting of the goal state and the other is iteratively separated into several states. The method is applied to a soccer robot which tries to shoot a ball into a goal. Fig 2 (b) shows the results in which the final state space is projected into two dimensional space in terms of the ball size and the goal size (when their positions are frontal and the orientation of the goal is horizontal). The grid lines indicate state segmentation designed by the programmer that are quite different in shape and size from the obtained states.

### 2.2 An On-line Learning Method

The above method needs sufficient amount of uniformly sampled data to construct the state space suitable for the robot to perform the given task, and therefore, does not cope with dynamical changes happened in the environment. These problems are resolved by the second approach [3] which obtains a purposive behavior within less learning time by incrementally segmenting the sensor space based on the experiences of the robot. The incremental segmentation is performed by constructing local models in the state space, which is based on the function approximation of the sensor outputs to reduce the



(a) basic idea



(b) 2-D projection of the result of state space construction

**Fig. 2. A basic idea of state-action space construction and result**

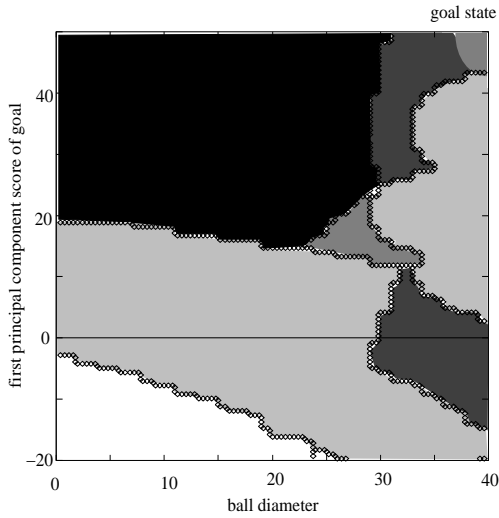
learning time, and the reinforcement signal to emerge a purposive behavior. They applied their method to the same task as in [2]. The basic ideas are as follows:

1. Set up a state space consists of two states; the goal state and the other.
2. Apply function approximation to the changes of the input vectors caused by action primitives. If the function approximation cannot cope with these changes, then segment the states into two and apply the function approximation to a new state. This process might cause to merge a state with one of the separated states. These processes can reduce ineffective explorations.
3. Initialize the action-value for the new state, and apply the reinforcement learning. The learning time is very short because the number of states to be updated is small.
4. Apply stochastic action selection to cope with dynamic change of the environment.

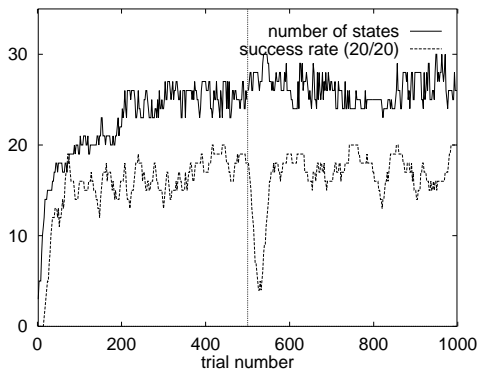
Figures 3 (a) and (b) show the experimental results. Figure 3 (a) shows a projection of the state space after 1,110 trials, where the state space in term of ball size and goal size is indicated when the position of the ball and the goal are center of the screen and the orientation of the goal is frontal. As we can see the shape of the resultant state spaces complicated and quite different from the previous result (see Figure 2 (b)). Figure 3 (b) indicates the changes of the success rate and the number of states in the case that the ball size is suddenly changed twice at the 500th trial. These suggest that the method cope with non-linear mapping between states and actions and deal with dynamic change of the environment.

### 3. A Scaling-Up Problem

Since each species of animals can be regarded to have its own intelligence, difference of intelligence seems to depend on the agent (capabilities in sensing, acting, and cognition) and its environment. If agents have the same bodies, differences or levels in intelligence can occur in the complexity of interactions with their environments. In case of our soccer playing robot with vision, the complexity of interactions may change due to other agents in the field such as common side players, opponents, judges and so on. In the following, we attempt at showing our view about the levels of complexity of interactions, especially from a viewpoint of the existence of other agents.



(a) state space



(b) success rate and the number of states

**Fig. 3. Experimental results**

1. **Self body and Static Environment:** The self body or static environment can be defined in a sense that the observable parts of which changes in the image plane can be directly correlated with the self motor commands (ex. looking at your hand showing voluntary motion, or observing an optical flow of the environment when changing your gaze). Theoretically, discrimination between “self body” and “static environment” is a hard problem because the definition of “static” is relative and depends on the selection of the reference (the base coordinate system) which also depends on the context of the given task. Usually, we suppose the natural orientation of the gravity and therefore it provides the ground coordinate system.

2. **Passive agents:** As a result of actions of the self or other agents, passive agents can be moving or stopped. A ball is a typical one. As long as they are stationary, they can be categorized into the static environment. But, not so simple correlation with motor commands as the self body or the static environment can be expected when they are in motion.

3. **Active (other) agents:** Active other agents do not have a simple and straightforward relationship with the self motions. In the early stage, they are treated as noise or disturbance because of not having direct visual correlation with the self motor commands. Later, they can be found as having more complicated and higher correlation (coordination, competition, and others). The complexity is drastically increased.

According to the complexity of the environment, the internal structure of the robot should be higher and more complex to emerge various intelligent behaviors. We show one of such structure coping with the complexity of agent-environment interactions with real robot experiments and discuss the future issues.

### 3.1 A More Complicated Task in Multi-Agent Environment

In a multi-agent environment, the conventional reinforcement learning algorithm does not seem applicable because the learner’s sensory information may change regardless of the learner’s motion due to the motion of other active agents in the environment. Therefore, the learner cannot predict the other agent behaviors correctly unless explicit communication is available. It is important for the learner to discriminate the strategies of the other agents and to predict their movements in advance to learn the behaviors successfully.

The existing methods in multi agent environments (ex., [8],[9],[10],[11],[12] and so on.) need state vectors in order for the learning to converge. However, it is difficult to obtain a reasonable analytical model in advance. Therefore, the modeling architecture is required to make the reinforcement learning applicable.

Here, we show a method which estimates the relationship between the learner’s behaviors and the other agents through interactions (observation and action) using the method of system identification. In order to construct the local pre-

dictive model of other agents, we apply Akaike’s Information Criterion(AIC) [13] to the result of Canonical Variate Analysis(CVA) [14], which is widely used in the field of system identification. The local predictive model is constructed based on the observation and action of the learner (observer).

We apply the proposed method to a simple soccer-like game in a physical environment. The task of the agent is to shoot a ball which is passed back from the other agent. Since the environment consists of the stationary agents (the goal and the line), a passive agent (the ball) and an active agent (the passer), the learner has to construct the adequate models for these agents. After constructing the models and estimating their parameters, the reinforcement learning is applied in order to acquire purposive behaviors. The proposed method can cope with the moving ball because a state vector for learning is selected appropriately so as to predict the successive steps.

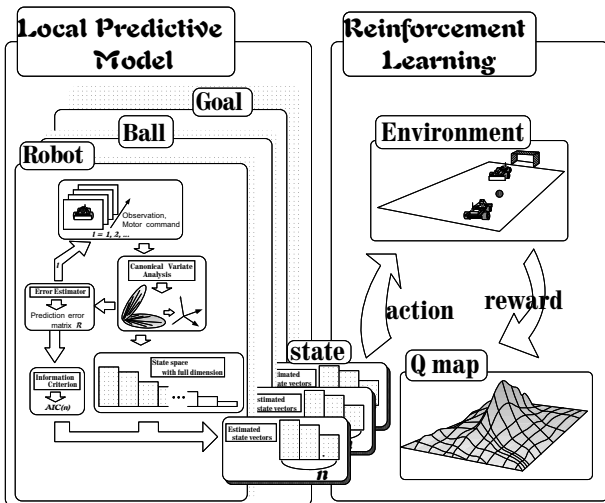


Fig. 4. An overview of the proposed method

Figure 4 shows an overview of the proposed method consisting of local predictive models and reinforcement learning architecture. At first, the learning agent collects the sequence of sensor outputs and motor commands to construct the local predictive models. By approximating the relationship between inputs (learner’s action) and outputs (observation), the local predictive model gives the learning agent not only the successive states of the agent but also the priority of state vectors, which means that first a few vectors might be sufficient to predict the successive states.

The flow of the proposed method is summa-

rized as follows:

1. Collect the observation vectors and the motor commands.
2. Estimate the state space with the full dimension directly from the observations and motor commands (Section 3.1.1).
3. Determine the dimension of the state vectors which is the result of the trade off between the error and the complexity of the model.
4. Apply the reinforcement learning based on the estimated state vectors.

### 3.1.1 Canonical Variate Analysis

A number of algorithms to identify multi-input multi-output (MIMO) combined deterministic-stochastic systems have been proposed [15]. In contrast to ‘classical’ algorithms such as PEM (Prediction Error Method), the subspace algorithms do not suffer from the problems caused by a priori parameterizations. Larimore’s Canonical Variate Analysis (CVA) [14] is one of such algorithms, which uses canonical correlation analysis to construct a state estimator (The details of the method is described elsewhere [16]).

### 3.1.2 Experimental Results

The output (observed) vectors are shown in Figure 5. In case of the ball, the center posi-

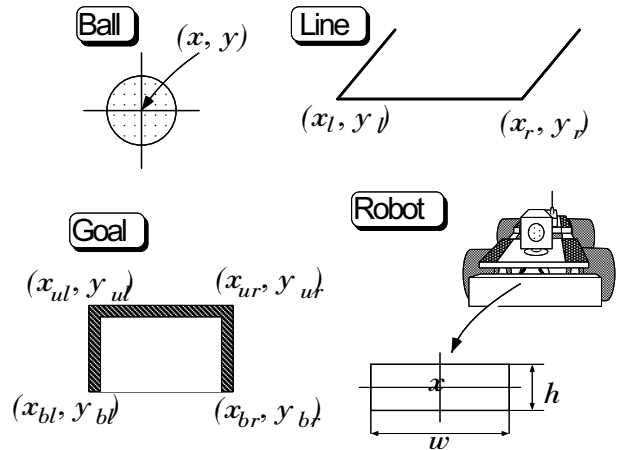


Fig. 5. Image feature points of the ball, goal, line and agent

tion of the ball image  $(x_c, y_c)$  is used, and the both ends  $(x_l, y_l)$  and  $(x_r, y_r)$  are used for the field lines. In case of the goal, the four corners of the goal image  $(x_{ul}, y_{ul})$ ,  $(x_{bl}, y_{bl})$ ,  $(x_{ur}, y_{ur})$ , and  $(x_{br}, y_{br})$  are used. In case of other agent, the center of position, the width and the height

of the plate are used. As a result, the dimension of the observed vector for the ball, the goal, the line, and the agent are 2, 4, 8, and 3 respectively.



Fig. 6. The real experiment set up with a passer and a shooter

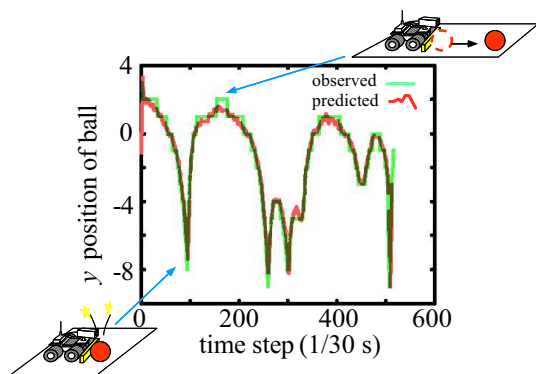


Fig. 7. Observed and predicted trajectories of the  $y$  position of the ball

Figures 6 and 7 show the real experiment set up with a passer and a shooter, and the observed and predicted trajectories of the  $y$  position of the ball, respectively. The ball trajectory is almost correctly predicted even though the learner kicked the ball at the time step 100.

#### 4. Concluding Remarks

Along with examples of soccer robots, we have claimed the importance of the design of the internal structure which reflects the complexity of the interactions with the agent's environment. Although the task and the environment seem simple and limited, the design of the soccer robots includes a variety of the fundamental and important issues as a standard problem in robotics and AI [1]. We expect that more agents in the field cause much higher interactions among them, which emerges a variety of more complex behaviors.

#### Acknowledgment

We like to thank Eiji Uchibe, Yasutake Takahashi, Masateru Nakamura, and Chizuko Mishima for their supports of our work described in this paper.

#### References

- [1] H. Kitano, M. Asada, Y. Kuniyoshi, I. Noda, E. Osawa, and H. Matsubara. "robocup: A challenge problem of ai". *AI Magazine*, 18:73–85, 1997.
- [2] M. Asada, S. Noda, and K. Hosoda. Action-based sensor space categorization for robot learning. In *Proc. of IEEE/RSJ International Conference on Intelligent Robots and Systems 1996 (IROS '96)*, pages 1502–1509, 1996.
- [3] Y. Takahashi, M. Asada, and K. Hosoda. Reasonable performance in less learning time by real robot based on incremental state space segmentation. In *Proc. of IEEE/RSJ International Conference on Intelligent Robots and Systems 1996 (IROS96)*, pages 1518–1524, 1996.
- [4] E. Uchibe, M. Asada, and K. Hosoda. Vision based state space construction for learning mobile robots in multi agent environments. In *Proceedings of 6-th European Workshop on Learning Robots, EWLR-6*, pages 33–41, 1997.
- [5] C. J. C. H. Watkins and P. Dayan. "Technical note: Q-learning". *Machine Learning*, 8:279–292, 1992.
- [6] R. S. Sutton. "Special issue on reinforcement learning". In R. S. Sutton (Guest), editor, *Machine Learning*, volume 8, pages –. Kluwer Academic Publishers, 1992.
- [7] M. Mataric. "Reward functions for accelerated learning". In *Proc. of Conf. on Machine Learning-1994*, pages 181–189, 1994.
- [8] Peter Stone and Manuela Veloso. Using machine learning in the soccer server. In *Proc. of IROS-96 Workshop on Robocup*, 1996.
- [9] E. Uchibe, M. Asada, and K. Hosoda. Behavior coordination for a mobile robot using modular reinforcement learning. In *Proc. of IEEE/RSJ International Conference on Intelligent Robots and Systems 1996 (IROS96)*, pages 1329–1336, 1996.
- [10] Michael L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *Proc. of Conf. on Machine Learning*.

*Learning-1994*, pages 157–163, 1994.

- [11] Tuomas W. Sandholm and Robert H. Crites. On multiagent Q-learning in a semi-competitive domain. In *Workshop Notes of Adaptation and Learning in Multiagent Systems Workshop, IJCAI-95*, 1995.
- [12] Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine Learning*, 8:293–321, 1992.
- [13] H. Akaike. A new look on the statistical model identification. *IEEE Trans. AC-19*, pages 716–723, 1974.
- [14] W. E. Larimore. Canonical variate analysis in identification, filtering, and adaptive control. In *Proc. 29th IEEE Conference on Decision and Control*, pages 596–604, Honolulu, Hawaii, December 1990.
- [15] Peter Van Overschee and Bart De Moor. A unifying theorem for three subspace system identification algorithms. *Automatica*, 31(12):1853–1864, 1995.
- [16] E. Uchibe, M. Asada, and K. Hosoda. “state space construction for behavior acquisition in multi agent environments with vision and action”. In *Proc. of ICCV 98*, pp.570–575, 1998.