

# Hierarchical Modular Reinforcement Learning Method and Knowledge Acquisition of State-Action Rule for Multi-target Problem

Takumi Ichimura

Faculty of Management and Information Systems,  
Prefectural University of Hiroshima  
1-1-71, Ujina-Higashi, Minami-ku,  
Hiroshima, 734-8559, Japan  
Email: ichimura@pu-hiroshima.ac.jp

Daisuke Igaue\*

Graduate School of Comprehensive Scientific Research,  
Prefectural University of Hiroshima  
\* He graduated from Prefectural Univ. of Hiroshima  
and is working at Iyo Bank, Ltd., Japan  
Email: punch20@gmail.com

**Abstract**—Hierarchical Modular Reinforcement Learning (HMRL), consists of 2 layered learning where Profit Sharing works to plan a prey position in the higher layer and Q-learning method trains the state-actions to the target in the lower layer. In this paper, we expanded HMRL to multi-target problem to take the distance between targets to the consideration. The function, called ‘AT field’, can estimate the interests for an agent according to the distance between 2 agents and the advantage/disadvantage of the other agent. Moreover, the knowledge related to state-action rules is extracted by C4.5. The action under the situation is decided by using the acquired knowledge. To verify the effectiveness of proposed method, some experimental results are reported.

**Index Terms**—Reinforcement Learning, Profit Sharing, Q-learning, Hierarchical Modular Reinforcement Learning, Multi-target, C4.5, Knowledge Acquisition

## I. INTRODUCTION

Multi-Agent Systems (MAS) where there a number of autonomous agents interacting with each affecting the actions of the other agents is a complex system. Learning enables MAS to be more flexible and robust and makes agents better able to handle uncertain and changing circumstances. Thus how to coordinate the behaviors of different agents by learning method is required. Reinforcement learning is an area of machine learning in computer intelligent system [3], [1], [2]. One of problems of reinforcement learning application of actual sized problem is “curse of dimensional problem”. High dimension of input leads to huge number of rules in the reinforcement learning application.

In order to solve these problems several types of hierarchical reinforcement learning have been proposed to apply actual applications [6], [7]. Hierarchical Modular Reinforcement Learning (HMRL), consists of 2 layered learning where Profit Sharing works to plan a prey position in the higher layer and Q-learning method trains the state-actions to the target in the

lower layer. In this paper, we expanded HMRL to multi-target problem under the consideration of the distance between targets. The function, called ‘AT field’, can estimate the interests for an agent according to the distance between 2 agents and the advantage/disadvantage of the other agent. Moreover, the knowledge related to state-action rules is extracted by C4.5. The action under the situation is decided by using the acquired knowledge. To verify the effectiveness of proposed method, some experimental results are reported.

The remainder of this paper is organized as follows. Section II describes about reinforcement learning method. Hierarchical modular reinforcement learning method is explained in Section III. In the section, we explain the multi-agent pursuit problem. Moreover, we give consideration to deal with the value of target according to the distance between 2 prey agents. Section IV is the knowledge discover of learning agents in the format of If-Then rules. In Section V, we give some discussions to conclude this paper.

## II. REINFORCEMENT LEARNING

The Profit Sharing and Q-Learning method are very popular in Reinforcement Learning. The section describes the algorithms of two kinds of Reinforcement Learning methods briefly.

### A. Profit Sharing

Multi agent systems have been developed in the field of Artificial Intelligence. Each agent is designed to work some schemes based on many rules which indicate knowledge of the agent world or relationship among the agents. However, the knowledge or relationship is not always effective to survive in their environment, because the agent will discard a partial of knowledge if its environment changes dynamically. Reinforcement Learning [3] is known to be worth to realize the cooperative behavior among agents even if little knowledge is provided with initial condition. The multi-agent system works to share a given reward among all agents.

Especially, PS method [1], [2] is an effective exploitation of reinforcement learning to adapt to a given environment.

©2013 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

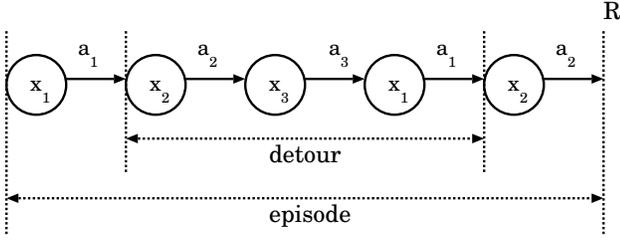


Fig. 1. The episode and the detour

In PS, an agent learns a policy based on the reward that is received from the environment when it reaches a goal state. It is important to design a reinforcement function that distributes the received reward to each action rule in the policy. In PS, the rule  $r_i$  is  $(s, a)$  for possible action  $a$  to a given sensory input  $x$  to  $s$ . The rule “If  $x$  then  $a$ .” is also written by  $\vec{x}\vec{a}$ . PS does not estimate the value function and computes weight of rules  $S_{r_i}$  for  $(s, a)$ . The episode is determined from the start state to the terminal state which the agent achieves the goal at time  $i$  and then a reward  $\mathbf{R}$  is provided. The PS gives the partial reward of  $\mathbf{R}$  to the fired rule  $(s_i, a_i)$  in an episode ( $i < W$ ).  $W$  is the maximum length of episode. The partial  $\mathbf{R}$  is determined by the value function  $f(i, \mathbf{R}, W)$ . Each rule is reinforced by the sum of current weight and slanted reward. That is,

$$S_{r_i} = S_{r_i} + f_i, \quad i = 0, 1, \dots, W - 1, \quad (1)$$

where  $S_{r_i}$  means the weight of the  $i$ th rule of an episode,  $f_i$  is the reinforce function and means the reinforce value at the  $-i$  step from obtaining  $\mathbf{R}$ .

The detour as shown in Fig.1 is the sequence of rules when the difference rules are selected for the same sensory input. There is a detour  $(\vec{x}_2\vec{a}_2, \vec{x}_3\vec{a}_3, \vec{x}_1\vec{a}_1)$  in the sequence  $(\vec{x}_1\vec{a}_1, \vec{x}_2\vec{a}_2, \vec{x}_3\vec{a}_3, \vec{x}_1\vec{a}_1, \vec{x}_2\vec{a}_2)$  in Fig.1. The rules in the detour may occur some ineffective rules. The ineffective rule is always on the detour from the episode. The other rules are called the effective rule. If the competition between ineffective rules and effective rules exists, the ineffectiveness are not reinforced. If the reinforcement function satisfies the ineffective rule suppression theorem, the reinforcement function is able to distribute more reward to effective rules than ineffective ones. In order to suppress such ineffective rules, the forgettable PS method is proposed.

$$L \sum_{j=1}^w f_j < f_{i-1}, \quad \forall i = 1, 2, \dots, W, \quad (2)$$

where  $f_i$  is the reinforcement function and  $L$  is the maximum number of effective rules. The reinforcement function decreases in a geometric series in the following.

$$f_i = \frac{1}{M} f_{i-1}, \quad i = 1, 2, \dots, W - 1, \quad (3)$$

where  $M (\geq L + 1)$  is a discount rate. Eq.(3) reinforces the rule from  $i = 1$  to  $i = W$  in an episode. Eq.(3) satisfied with the curve as shown in Fig. 3.

The algorithm of PS is as follows.

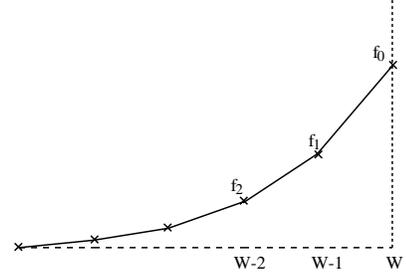


Fig. 3. Reinforcement Function

- Step 1) Initialize  $S_{r_i}$  arbitrarily.

Step 2) Repeat (for each episode):

  - a) Initialize  $r_i$  and  $W$ .
  - b) Repeat (for each step of episode):
    - i)  $a \leftarrow$  action given by  $\pi$  for  $\mathcal{S}$  at state  $x$
    - ii) Take action  $a$ ; observe reward,  $\mathbf{R}$  and next state  $\hat{x}$
    - iii)  $\forall i, i \leftarrow i + 1$ , set  $r_0 = \vec{x}\vec{a}$
    - iv) If  $R \neq 0$ , set  $f_0 = \mathbf{R}$  and calculate the following.
$$S_{r_i} = S_{r_i} + f_i, \quad i = 0, 1, \dots, W - 1, \quad (4)$$

where  $f_i = \frac{1}{M} f_{i-1}, i = 1, 2, \dots, W - 1.$
    - v)  $x \leftarrow \hat{x}, W = W + 1$
  - c) until  $x$  is terminal

Fig. 2. The algorithm of PS

## B. Q-Learning

Temporal Difference (TD) method can directly learn from raw experience without a model of the environment’s dynamics [3]. TD method uses experience to solve the prediction problem. If a non-terminal state  $s_t$  is visited at time  $t$ , TD method updates their estimate  $V(s_t)$  based on events after that visit. TD method waits only until the next time step. That is, TD method forms a target at time  $t + 1$  and makes an appropriate update using the observed reward  $r_{t+1}$  and the estimate  $V(s_{t+1})$ . The simplest expression in TD method can be written as follows.

$$V(s_t) \leftarrow V(s_t) + \alpha [r_{t+1} + \gamma V(s_{t+1}) - V(s_t)] \quad (5)$$

TD method can learn their estimates in part on the basis of other estimates. TD method is used for the evaluation or prediction by applying generalized policy iteration. We use the Q-learning as an off-policy TD method in this paper, because the learned action-value function,  $Q$ , directly approximates the optimal action-value function,  $Q^*$  with no dependence of policy. The simplest Q-learning can be written as follows.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)] \quad (6)$$

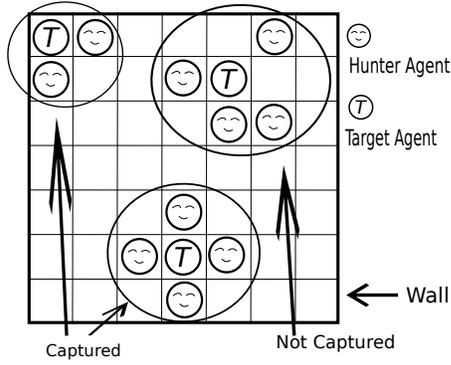


Fig. 5. Multi-Agent Pursuit Problem

- Step 1) Initialize  $Q(s, a)$  arbitrarily.  
 Step 2) Repeat (for each episode):
- a) Initialize  $s$ .
  - b) Repeat (for each step of episode):
    - i) Choose  $a$  from  $s$  by using policy derived from  $Q$
    - ii) Take action  $a$ ; observe  $r$  and next state  $\acute{s}$
    - iii) Eq.(6) is executed.
    - iv)  $s \leftarrow \acute{s}$ ;
  - c) until  $x$  is terminal

Fig. 4. The algorithm of Q-Learning

### III. HIERARCHICAL MODULAR REINFORCEMENT LEARNING METHOD

This section defines Multi-Agent Pursuit Problem to explain the simulation environment where the Hierarchical Modular Reinforcement Learning (HMRL) [7] Method works. Moreover, we develop the HMRL method to work in Multi-Agent Pursuit Problem where two or more kinds of prey agents works in the same environment.

#### A. Multi-Agent Pursuit Problem

The pursuit problem is well-known to be an appropriate example of cooperative Mult-agent system (MAS) [5]. In this study, the pursuit problem is considered in a  $7 \times 7$  grid world, where two prey agents( $T$ ) and four hunter agents are placed at random positions in the environment as shown in Fig.5. Hunters are learning agents and try to capture the randomly moving prey. In this paper, the prey agent does not learn the state-action rule through the experience and the two or more prey agents does not work to cooperate with each other. At each time, agents synchronously select and perform on out of five actions without communicating with each other: Staying at the current position or moving north, south, west, or east. Preys and hunters cannot share a cell. Also, an agent is not allowed to move off the environment. The prey is captured, when all of its neighbor cells are occupied by hunters as shown in Fig.5.

#### B. Hierarchical Modular Reinforcement Learning Method

For the pursuit problem, huge memory consumption is required to express the internal knowledge of the agents.

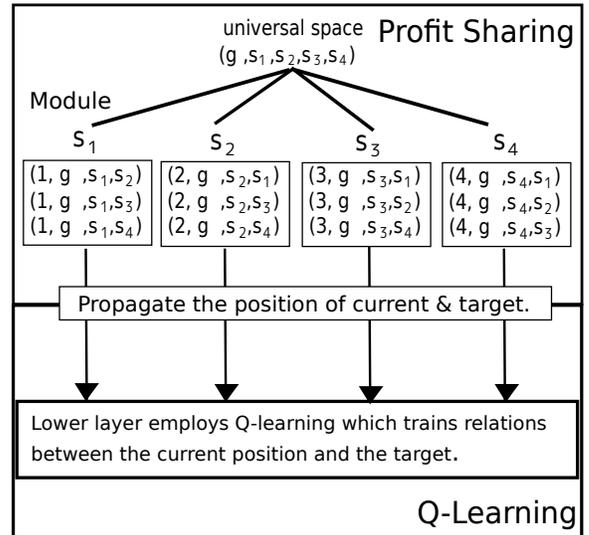


Fig. 6. Hierarchical Reinforcement Learning

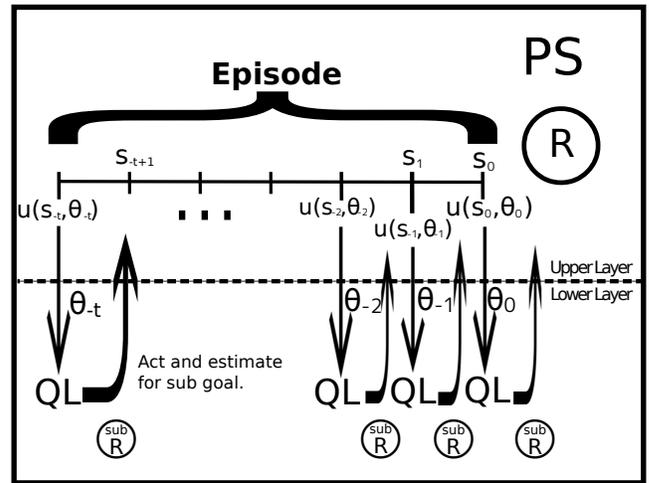


Fig. 7. A Distribution of Reward in Hierarchical Reinforcement Learning

Moreover, because the surrounding environment is complex, the agents cannot express the collaboration. [6], [7] proposed the hierarchical modular reinforcement learning to solve the above problems. It is difficult to decide how many kinds of sub-task should be decomposed into.

In [7], Prof. Watanabe conceived of the idea that decomposes the surrounding task(capturing) into “decision of move position target” for surrounding according to current monitored state and “selection of appropriate action” to move to the target position of each agent. The task is decomposed into “surrounding” task synchronized with the other hunter agents and “exploring the environment” task. Moreover, the upper task corresponds only to collaborative surrounding strategy.

In the upper layer, the target position of the agent is decided based on observed state such as the current position of the prey agent and the other hunter agents. The rules in the upper layer express goodness of the target position corresponding to the current state excluding actual actions. In order to construct the rules based on the current state combination, huge

corresponding memory is needed. To avoid such requirement, the authors applied modular structure for the rule expression [7] in the upper layer as shown in Fig.6. In Fig.6, the state space is divided to 4 sub-spaces where the following equation is satisfied.

$$(g, s_1, s_2, s_3, s_4) = \cup_e (e, g, s_e, s_\epsilon), (e, \epsilon \in E, e \neq \epsilon) \quad (7)$$

The weights of rules in the upper layer are updated by Profit Sharing as follows.

$$\begin{aligned} u(e, g(i), h_e(i), h_\epsilon(i)) &= u(e, g(i), h_e(i), h_\epsilon(i)) \\ &+ k(e, g(i), h_e(i), h_\epsilon(i)), \\ k(e, g(i-1), h_e(i-1), h_\epsilon(i-1)) &= \\ \rho k(e, g(i), h_e(i), h_\epsilon(i)) & \end{aligned} \quad (8)$$

$(i = 0, -1, \dots, -m, \epsilon \neq e),$

where  $u(\cdot)$  is the estimate function for target position and  $k(\cdot)$  is an reinforcement function as shown Fig.3.  $e$  is the hunter agent and  $\epsilon$  is the other hunter agent.  $g(i)$  is the position and  $h_e(i)$  is the position at time  $i$ , respectively. Time 0 is when the hunter agent receives the reward.  $\rho$  is the parameter.

The target position is divided as a sub goal for surrounding tasks instead of final goal corresponding to the current state of the prey agent according to the weight of rules. The target position of the agent is determined by the following equation.

$$\theta_e = \arg \max_v \sum_{\epsilon} \frac{u(e, g, v, h_\epsilon)}{\mu^{|h_e - v|}}, (\epsilon \neq e, \mu \geq 1),$$

where  $v$  is the candidate of target position. According to the selected state, the information for target position is sent to the lower layer.

In the lower layer, the selection of action to walk to the target position decided at the upper layer is implemented by reinforcement learning process as Q-learning:

$$\begin{aligned} Q(s_e(t), a_e(t), \theta_e) &= Q(s_e(t), a_e(t), \theta_e) \\ &+ k(r_t + \gamma \max_{\eta} Q(s_e(t+1), \eta, \theta_e) - Q(s_e(t), a_e(t), \theta_e)), \end{aligned}$$

where  $Q$  is  $Q$ -value,  $s_e(t)$  and  $a_e(t)$  are the state vector and the action of the agent  $e$  at  $t$ th step, respectively.  $\theta_e$  is the position of agent  $e$ .  $r_t$  is the reward.  $\eta$  is the maximum value of action.  $k$  is the step size parameter.

### C. 2 prey agent based Hierarchical Reinforcement Learning

Multi-Agent Pursuit Problem with 2 prey agents in the environment is discussed in this paper. For the problem, we consider the division of space as shown in Fig.8. If there are 2 prey agents, the environment has 2 goals. Therefore, the relation among sub spaces in Fig.8 is defined as follows:

$$\begin{aligned} (g_1, g_2, s_1, s_2, s_3, s_4) &= \cup_e \cup_l (e, g_l, s_e, s_\epsilon) \\ (e, \epsilon \in E, l \in L, e \neq \epsilon), & \end{aligned} \quad (9)$$

where  $g_l$  is the goal position for each prey agent.

Each modular has 2 target position, but only one target position should be sent to the lower layer. Therefore, the judgment rule for the decision of appropriate position is defined as follows:

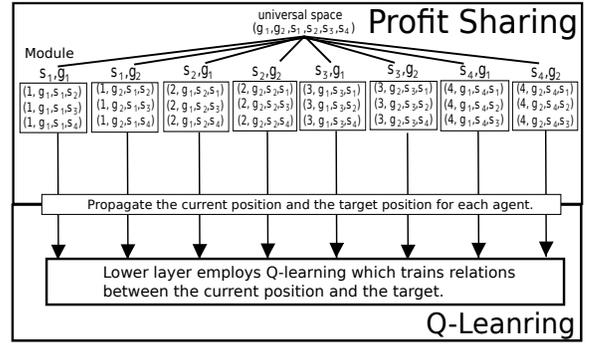


Fig. 8. 2 Prey Agent Model

$$\theta_e = \begin{cases} \arg \max_v \sum_{\epsilon} \frac{u(e, g_0, v, h_\epsilon)}{\mu^{|h_e - v|}} & \text{if } |h_e - g_0| < |h_e - g_1| \\ \arg \max_v \sum_{\epsilon} \frac{u(e, g_1, v, h_\epsilon)}{\mu^{|h_e - v|}} & \text{if } |h_e - g_1| < |h_e - g_0| \end{cases} \quad (\epsilon \neq e, \mu \geq 1) \quad (10)$$

If the target is quite different altering behaviors among the prey agents, e.g. a target has positive reinforcement and the other agent has punishment, it is difficult to consider the value of target simultaneously. In this paper, Eq.(11) is defined by the idea that when there are 2 kinds of target due to their value, positive and negative, the value of  $u(\cdot)$  is changed according to the degree to be affected by each other.

$$ATF = \begin{cases} \Phi = 0.0 & \text{(if } gd \leq n_1) \\ \Phi = 1.0 & \text{(if } n_1 < gd \leq n_2) \\ \Phi = 0.9 & \text{(if } n_2 < gd) \end{cases}, \quad (11)$$

where  $gd$  is the distance between two agents.  $n_1$  is the parameter to judge for whether the distance of the agent and the other agent is within close distance and  $n_2$  is the parameter to judge whether the distance is long distance. In this paper, we set that  $n_1 = 2$  and  $n_2 = 5$ . The estimate value is updated by using Eq.(12).

$$\begin{aligned} u(e, g_l(i), h_e(i), h_\epsilon(i)) &= u(e, g_l(i), h_e(i), h_\epsilon(i)) \\ &+ k(e, g_l(i), h_e(i), h_\epsilon(i)) \\ k(e, g_l(i-1), h_e(i-1), h_\epsilon(i-1)) &= \rho \cdot ATF(gd) \cdot k(e, g_l(i), h_e(i), h_\epsilon(i)) \\ (e, \epsilon \in E, l \in L, i = 0, -1, \dots, -m), & \end{aligned} \quad (12)$$

where  $ATF(gd)$  is the function of AT-Field given by Eq.(11). The output of function can be reduced by the discount factor,  $\rho$ , according to the degree that the corresponding agent is affected by the other agent.  $e$  and  $\epsilon$  are the index of hunter and prey agent, respectively.  $E$  and  $L$  mean the set of all agents and the prey agent, respectively.  $h$  and  $g$  are the position of the hunter agent and the prey agent, respectively.  $ATF$  function will not affect the division of state space in the profit sharing.

### D. Simulation Results

This section describes the simulation results under the 2 prey agent and 4 hunter agents. The position of all agents are

TABLE I  
STEP NUMBER WITHOUT ATF

Iterations	Episode		Action	
	Ave.	Var.	Ave.	Var.
1-200	684.7	2374.0	958.6	1596.7
201-2,000	355.2	127.9	358.5	134.5
2,001-17,000	101.3	5.1	104.9	5.3
17,001-20,000	62.3	3.3	66.6	3.1

randomly assigned in the 7 grid. A trial is starting from the initial situation until the hunter agents capture 2 prey agents as shown in Fig.5. After one trial the environment and  $Q$ -value are initialized, and a set of simulation is till 20,000 trials. The reward is 100 if the prey agent is positive target and it is 0 otherwise. In lower layer, when the agent reaches to the target position sent from the upper layer, the agent can receive the reward 100. The behavior of prey agent is randomly and the hunter agent moves due to the acquired state-action rules. Of course, each agent does not know the behaviors of the other agents.

In order to evaluate the effectiveness of the proposed model, we define the 3 ratio of capturing targets: “(Within Safety)”, “(Within Dangerous)”, and “(Positive Ratio)”.

- 1) (Within Safety): When the hunter agent captured a prey agent with positive reward, if the distance between them is larger than  $n_1$ , the captured prey agents belongs to the set  $far$ .

$$P(\text{safety\_distance}) = \frac{\#(\text{safety\_target} \cap \text{far})}{\#(\text{safety\_target})} \quad (13)$$

- 2) (Within Dangerous): When the hunter agent captures a prey agent with positive reward, if the distance between them is smaller than  $n_1$ , the captured prey agents belongs to the set  $near$ .

$$P(\text{dangerous\_distance}) = \frac{\#(\text{safety\_target} \cap \text{near})}{\#(\text{safety\_target})} \quad (14)$$

- 3) (Positive Ratio): It means the ratio of (Within Safety) over the simulations.

$$P(\text{safety\_target}_{\text{positive}}) = \frac{\#(\text{safety\_target} \cap \text{far})}{\text{iteration}} \quad (15)$$

Table I and Table II show that the number of steps and the actions without ATField model and with one, respectively, until the prey target is captured. The simulation result related to the number of steps and actions are almost same results, although the computation time with ATField model gets longer than that without ATField model.

Table III and Table IV show that the capture ratio of (Safety Target), (Within Safety), (Within Dangerous), (Positive Ratio), and distance, without ATField model and with one, respectively. The distance means the distance between targets. From these tables, the performance in model with ATField is better than that without ATField model.

#### IV. KNOWLEDGE ACQUISITION

The state-action rules in the lower layer are extracted by C4.5. Fig.9 shows the part of extracted results. The rules

TABLE II  
STEP NUMBER WITH ATF

Iterations	Episode		Action	
	Ave.	Var.	Ave.	Var.
1-200	703.5	2881.1	999.1	2372.9
201-2,000	403.7	193.4	406.7	203.7
2,001-17,000	119.7	4.2	122.7	4.6
17,001-20,000	74.5	3.7	77.7	3.6

are extracted while training the module. Fig.9 is the result of 19,900-20,000 trials. ‘theta\_x’ and ‘theta\_y’ mean the difference in the x axis and y axis in the move, respectively. The output in the teaching signal is the target position sent from the upper layer. The simulation is 72,327 instances in the teach data set.

```

54pt
theta_Y > -1
| theta_X <= -1
| | theta_Y <= 0: left (12519.0/1907.0)
| | theta_Y > 0
| | | theta_Y <= 1: left (2172.0/1096.0)
| | | theta_Y > 1
| | | | theta_X <= -2
| | | | | theta_Y <= 2: down (270.0/128.0)
| | | | | theta_Y > 2
| | | | | | theta_X <= -3
| | | | | | | theta_X <= -5
| | | | | | | | theta_Y <= 3: left (4.0/1.0)
| | | | | | | | theta_Y > 3: down (2.0/1.0)
| | | | | | | | theta_X > -5
| | | | | | | | theta_Y <= 3: down (20.0/8.0)
| | | | | | | | theta_Y > 3: stay (5.0/2.0)
| | | | | | | | theta_X > -3: left (56.0/29.0)
| | | | | | | | theta_X > -2: down (648.0/328.0)
| theta_X > -1
| | theta_Y <= 0
| | | theta_X <= 0: stay (8056.0)
| | | theta_X > 0
| | | | theta_X <= 2: right (12260.0/1733.0)
| | | | theta_X > 2
| | | | | theta_X <= 3: right (320.0/179.0)
| | | | | theta_X > 3
| | | | | | theta_X <= 4: stay (442.0/104.0)
| | | | | | theta_X > 4: right (47.0/33.0)
| | | | | | | theta_Y > 0
| | | | | | | | theta_X <= 0
| | | | | | | | theta_Y <= 1: down (11541.0/1451.0)
| | | | | | | | theta_Y > 1
| | | | | | | | theta_Y <= 2: down (959.0/352.0)
| | | | | | | | theta_Y > 2
| | | | | | | | | theta_Y <= 3: down (328.0/199.0)
| | | | | | | | | theta_Y > 3
| | | | | | | | | | theta_Y <= 5: stay (173.0/83.0)
| | | | | | | | | | theta_Y > 5: left (11.0/4.0)

```

Fig. 9. Calculation Results of C4.5 (partial)

For easy comprehension, Fig.10 shows the extracted knowledge as shown in Fig.9 in the If-Then rule format. In the simulation, we can get 47 state-action rules. Fig.10 shows 10 sample rules only.

```

No.1
If theta_X <= 4 theta_X > 2 theta_Y <= -6 Then up
with CF=1.0

```

TABLE III  
CAPTURE RATIO OF PREY AGENTS WITHOUT ATF

	Safety Target		Within Safety		Within Dangerous		Positive ratio		Distance	
	Ave.	Var.	Ave.	Var.	Ave.	Var.	Ave.	Var.	Ave.	Var.
1-200	53.1%	14.4	60.2%	24.0	39.8%	24.0	32.0%	14.0	3.22	0.030
201-2000	74.8%	1.4	77.2%	2.7	22.8%	2.7	57.8%	4.0	3.95	0.004
2001-17000	86.4%	0.2	84.3%	0.1	15.7%	0.1	72.9%	0.3	4.27	0.001
17001-20000	84.7%	0.7	83.3%	0.4	16.7%	0.4	70.6%	0.9	4.11	0.001

TABLE IV  
CAPTURE RATIO OF PREY AGENTS WITH ATF

	Safety Target		Within Safety		Within Dangerous		Positive ratio		Distance	
	Ave.	Var.	Ave.	Var.	Ave.	Var.	Ave.	Var.	Ave.	Var.
1-200	53.7%	24.5	49.2%	13.9	50.8%	13.9	26.5%	16.6	2.93	0.03
201-2000	73.4%	1.2	77.2%	2.7	22.8%	2.7	56.7%	2.7	3.94	0.004
2001-17000	89.7%	0.1	86.9%	0.04	13.1%	0.04	77.9%	0.1	4.51	0.0005
17001-20000	90.6%	0.4	86.4%	0.5	13.6%	0.5	78.3%	0.3	4.41	0.001

TABLE V  
STEP NUMBER

	Episode	Action
1-200	543.1	543.4
201-2,000	194.0	195.6
2,001-17,000	63.0	69.3
17,001-20,000	46.3	54.6

TABLE VI  
CAPTURE RATIO OF PREY AGENT

	Safety Target	Within Safety	Within Dangerous	Positive ratio	Target Distance
1-200	56.5%	61.1%	38.9%	34.5%	3.27
201-2,000	85.7%	84.8%	15.2%	72.6%	4.44
2,001-17,000	92.2%	85.5%	14.6%	78.9%	4.36
17,001-20,000	91.4%	84.3%	15.7%	77.1%	4.23

```

No.2
If theta_X <= 0 theta_X > -1 theta_Y <= 0 theta_Y > -1 Then stay
with CF=1.0
No.3
If theta_X <= 0 theta_X > -1 theta_Y <= 1 theta_Y > 0 Then down
with CF=0.8742743263148774
No.4
If theta_X <= 2 theta_X > 0 theta_Y <= 0 theta_Y > -1 Then right
with CF=0.8586460032626427
No.5
If theta_X <= 0 theta_X > -1 theta_Y <= -1 Then up
with CF=0.8478816513050886
No.6
If theta_X <= -1 theta_Y <= 0 theta_Y > -1 Then left
with CF=0.8476715392603243
No.7
If theta_X <= 4 theta_X > 3 theta_Y <= 0 theta_Y > -1 Then stay
with CF=0.7647058823529411
No.8
If theta_X <= -5 theta_Y <= 3 theta_Y > 2 Then left CF=0.75
No.9
If theta_X <= 1 theta_X > 0 theta_Y <= 5 theta_Y > 4 Then stay
with CF=0.7272727272727273
No.10
If theta_X <= -6 theta_Y <= -1 theta_Y > -2 Then left
with CF=0.7142857142857143

```

Fig. 10. IF-Then rules (partial)

By using the acquired rules, the simulation results are the number of steps and the actions, and the ratio of captured prey agents as shown in Table V and VI. The performance with rules are better than that of without rules.

## V. CONCLUSIVE DISCUSSION

Hierarchical Modular Reinforcement Learning (HMRL)[7], consists of 2 layered learning where Profit Sharing works to plan a prey position in the higher layer and Q-learning method trains the state-actions to the target in the lower layer. If the

multi-agent pursuit problem has 2 or more prey agents, in many cases, the reward for them is set toward same purpose, that is, the rewards are same value. In this paper, we expanded HMRL to multi-target problem under the consideration of the distance between targets. The function, called 'AT field', can estimate the interests for an agent according to the distance between 2 agents and the advantage/disadvantage of the other agent. Moreover, the knowledge related to state-action rules is extracted by C4.5. In simulation results, AT field function is effective to measure the difference between the rewards of prey agents. We will verify the method in real world problem in future.

## REFERENCES

- [1] J.J.Grefenstette, *Credit Assignment in Rule Discovery Systems Based on Genetic Algorithms*, Machine Learning, Vol.3, pp.225-245, 1998.
- [2] K.Miyazaki, S.Arai, and S.Kobayashi, *A Theory of Profit Sharing in Multi-agent Reinforcement Learning*, Journal of Japanese Society for Artificial Intelligence, Vol.14, No.6, pp.1156-1164, 1999 (Japanese).
- [3] R.S.Sutton and G.B.Andrew, *Reinforcement Learning: An Introduction*, MIT Press, 1998.
- [4] Y.Ishiwaka, T.Sato, Y.Kakazu, *An approach to the pursuit problem on a heterogeneous multiagent system using reinforcement learning*, Robotics and Autonomous Systems, Vol.43, No.4, pp.245-256, 2003.
- [5] Z.Pu-Cheng, H.Bing-Rong, H.Qing-Cheng, and J.Khurshid, *Hybrid Multi-agent reinforcement Learning Approach:The Pursuit Problem*, Information Technology Journal, Vol.5, No.6, pp.1006-1011, 2006.
- [6] T.Wada, T.Okawa, T.Watanabe, *A study on hierarchical modular reinforcement learning for multi-agent pursuit problem based on relative coordinate states*, Proc. of the 8th IEEE international conference on Computational intelligence in robotics and automation, pp.302-308, 2009.
- [7] T.Watanabe and T.Wada, *A Study on Hierarchical Modular Reinforcement Learning Algorithm for Multi-Agent Pursuit Problem Based on Relative Coordinate States*, Journal of Bio-medical Fuzzy System Association, Vol.12, No.2, pp.65-74, 2010 (Japanese).