

**Region-based Q-learning
For Autonomous Mobile Robot Navigation**

(Tel : +82-31-408-5802; Fax : +82-31-408-5803;
E-mail: ihsuh@email.hanyang.ac.kr)

Abstract : Q-learning, based on discrete state and action space, is a most widely used reinforcement learning. However, this requires a lot of memory and much time for learning all actions of each state when it is applied to a real mobile robot navigation using continuous state and action space. Region-based Q-learning is a reinforcement learning method that estimates action values of real state by using triangular-type action distribution model and relationship with its neighboring state which was defined and learned before. This paper proposes a new Region-based Q-learning which uses a reward assigned only when the agent reached the target, and get out of the local optimal path with adjustment of random action rate. If this is applied to mobile robot navigation, less memory can be used and robot can move smoothly, and optimal solution can be learned fast. To show the validity of our method, computer simulations are illustrated.

Keywords : RQ-learning, neighboring state, action distribution model, mobile robot, navigation

1. Introduction

2.1 Q-learning

Q-learning is an off-policy reinforcement learning method. The optimal action $p^*(s_t)$ is defined as:

$$p^*(s_t) = \arg \max_a Q(s_t, a) \quad (1)$$

The Q-value is updated according to the Bellman optimality equation:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (r_{t+1} + g \max_a \{Q(s_{t+1}, a)\} - Q(s_t, a_t)) \quad (2)$$

where α is the learning rate, r_{t+1} is the reward, and g is the discount factor.

2.2 RQ-learning

RQ-learning is a region-based Q-learning method. It uses a triangular-type action distribution model and relationship with its neighboring state. The Q-value is updated according to the Bellman optimality equation:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (r_{t+1} + g \max_a \{Q(s_{t+1}, a)\} - Q(s_t, a_t))$$

where α is the learning rate, r_{t+1} is the reward, and g is the discount factor.

2. Q-learning RQ-learning

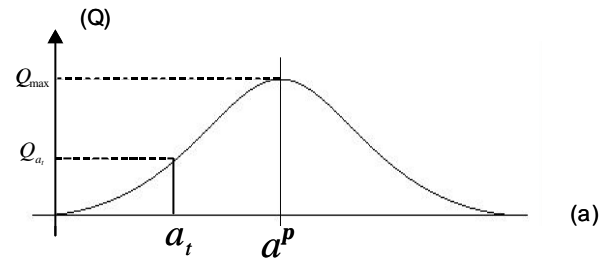
3. Region-based Q-learning (RQ-learning)

Q-learning

가

RQ-learning

()



2.

a_t

3.1

Q-learning

(action distribution model)

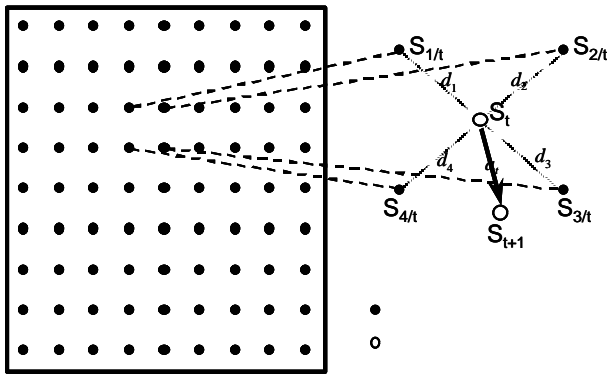
(Q_max)

(a^p)

가

가

(relationship)



1.

가

가

1

$$m_i = \frac{1-d_i}{\sum_j (1-d_j)}, \text{ if } (d_i > 1) \rightarrow d_i = 1 \quad (3)$$

3.2 (action distribution model)

Q-

가

가

가

가

Q-

distribution model)

가

가

(action

3.3 RQ-learning

[RQ-learning]

1.

2.

1)

2)

3)

4)

5)

6)

Q-learning

Q-

RQ-

learning

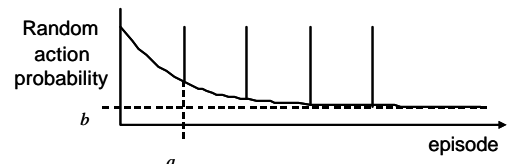
$$a_t = \sum_{i=1}^{2^N} m_i a_i, \text{ N:dimension} \quad (5)$$

Q-

가

가

가



3. episode

가

0

1

$$r = \begin{cases} 1 & \text{when it achieve the goal} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

. (4))

$$\max_a \{Q(s_{t+1}, a)\} = \max_a \left\{ \sum_{i=0}^{2^N} \mathbf{m}_i Q_i(s_{i/t+1}, a) \right\} \quad (7)$$

Watkin's action-

value update equation

$$Q_i(s_{i/t}, a_i) \leftarrow Q_i(s_{i/t}, a_i) + \mathbf{a} (r_{i/t+1} + \mathbf{g} \max_a \{Q(s', a)\} - Q(s_{i/t}, a_i))$$

$$r_{i/t+1} = \mathbf{m}_i r_{t+1}, \quad a_i : s_i \rightarrow s_{i+1}, \quad s_{i/t} \rightarrow s' \quad (8)$$

$\mathbf{a} = 1$

[1].

$$Q_i(s_{i/t}, a_i) \leftarrow r_{i/t+1} + \mathbf{g} \max_a \{Q(s', a)\} \quad (9)$$

4.

25*25

Q-learning RQ-learning

4

5

Q-learning

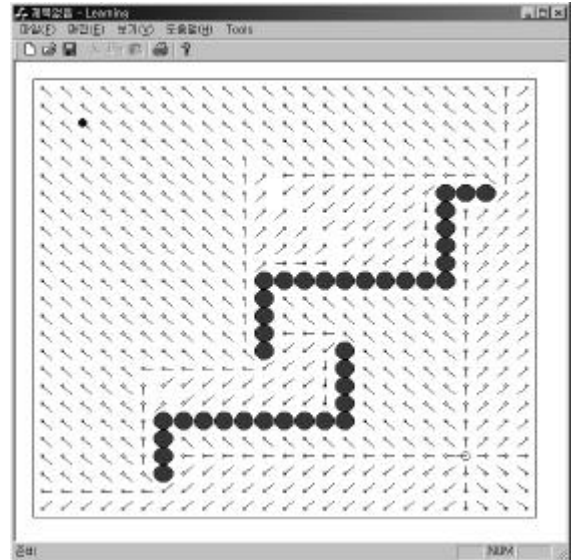
가

6

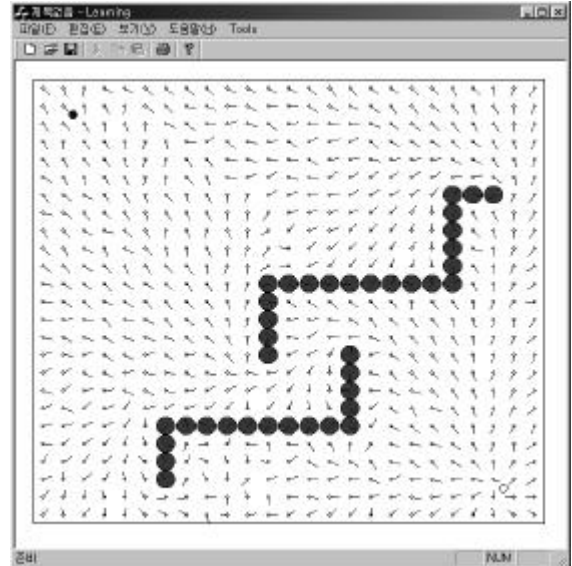
7

RQ-learning

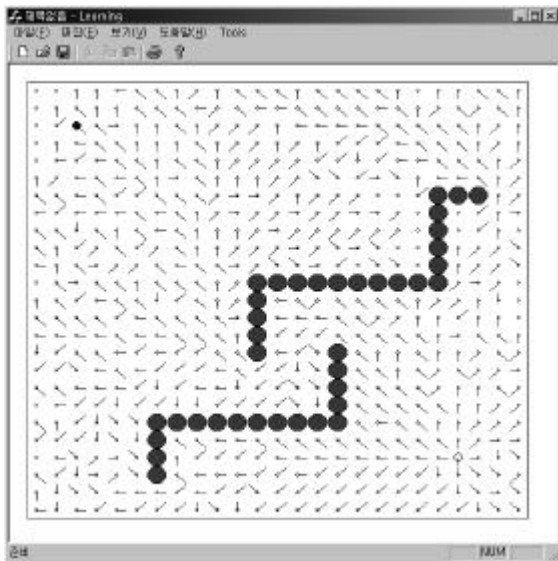
RQ-learning Q-learning



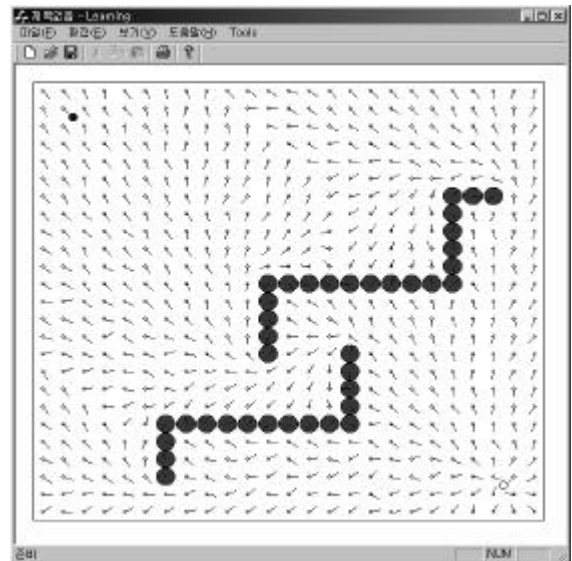
5. Q-learning , 124 episode



6. RQ-learning , 54 episode



4. Q-learning , 89 episode



7. RQ-learning , 62 episode

가 가

- 1.
- 2.

가

- 3.

가

- 1.

	Q-learning		RQ-learning	
	episode	step	episode	step
	89	57245	54	36214
	124	81969	62	48542

5.

Q-learning

RQ-learning

RQ-learning

Q-learning

RQ-learning

RQ-learning

가

Attolico, "Learning actions from vision-based positioning in goal-directed navigation," IROS'98 International Conference on Intelligent Robots and Systems, 1998.

[8] Yasutake Takahashi, Masanori Takeda, and Minoru Asada, "Continuous Valued Q-learning for Vision-Guided Behavior Acquisition," Proceeding of the 1999 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems, Taipei, Taiwan, R.O.C, pp. 2555-260, August, 1999.

[1] R. S. Sutton and A. G. Barto, Reinforcement Learning. Cambridge, MA : MIT Press, 1998.

[2] , " Q-learning," , pp. 271-276, 1997.

[3] C. Watkins, P. Dayan, "Q-learning, technical note," Machine Learning, Vol. 8, pp. 279-292, 1992.

[4] Yoichi Hirashima, "Q-learning Algorithm Using an Adaptive-Sized Q-table," Proceedings of the 38th Conference on Decision & Control, Phoenix, pp. 1599-1604, December, 1999.

[5] T. D'Orazio, G. Cicirelli, "Continuous Reward versus Discrete Reward in a Qlearning Agent," The Fifth ICARCV'98, Singapore, pp. 584-588, December, 1998.

[6] J. del, R. Millan, "Rapid, safe, and incremental learning of navigation strategies," Sys. Man and Cybernetics, 26(3), 1996.

[7] T. D'Orazio, G. Cicirelli, C. Distanto and G.