

[abtahi@aut.ac.ir](mailto:abtahi@aut.ac.ir) [mmeybodi@aut.ac.ir](mailto:mmeybodi@aut.ac.ir)

## **Solving Multi-Agent Markov Decision Processes Using Learning Automata**

**F. Abtahi      M. R. Meybodi**

Soft Computing Laboratory  
Department of Computer Engineering and Information Technology  
Amirkabir University of Technology, Tehran ,Iran

[abtahi@aut.ac.ir](mailto:abtahi@aut.ac.ir) [mmeybodi@aut.ac.ir](mailto:mmeybodi@aut.ac.ir)

**Abstract:** Multi-Agent Markov Decision Processes (MMDPs) are widely used to model many types of multi-agent systems. In this paper, several algorithms based on learning automata for solving MMDPs and finding a policy for selecting actions are proposed. In the proposed algorithms, Markov problem is described as a directed graph. The nodes of this graph are the states of the problem, and the directed edges represent the actions that result in transition from one state to another. Each node in the graph is equipped with a learning automaton whose actions are the outgoing edges of that node. Each agent moves from one node to another and tries to reach the goal state. In each node, the agent with the help of the learning automaton in that node chooses the next transition. The actions taken by the learning automata along the path traveled by the agent is then rewarded or penalized based on the cost of the traveled path according to a learning algorithm. This way the optimal policy for the agent will be gradually reached. The results of experiments have shown that the proposed algorithms perform better than the existing learning automata based algorithms in terms of cost and the speed of reaching the optimal policy.

**Keywords:** Multi- Agent Systems, Multi-Agent Markov Decision Process, Learning Automata, Optimal Policy

$$\begin{aligned}
& S = \{i\} \quad \langle S, A, P, R \rangle \quad (MDP) \quad [ ] \\
& \quad \quad \quad P: S \times A \times S \rightarrow [0,1] \quad A = \{a\} \\
& \quad \quad \quad R: S \rightarrow \mathbb{R} \quad P(i, a, j) \quad j \quad i \quad a \\
& \quad \quad \quad MDP \quad [ ] \quad R(i) \quad i \\
& \quad \quad \quad MDP \quad ( \quad ) \\
& \quad \quad \quad \pi \quad \pi: S \rightarrow A \\
& \quad \quad \quad M = \{m\} \quad n \\
& \quad \quad \quad MDP \quad A_m = \{a_m\} \quad S_m = \{i_m\} \\
& \quad \quad \quad MDP \quad S_M \quad A_M \\
& \quad \quad \quad \langle S_M, A_M, P_M, R_M \rangle \quad MDP \\
& \quad \quad \quad S_M = S_1 \times \dots \times S_n \quad A_M = A_1 \times \dots \times A_n \\
& \quad \quad \quad [ ] \quad R_M: S_M \rightarrow \mathbb{R} \quad P_M: S_M \times A_M \times S_M \rightarrow [0,1] \\
& \quad \quad \quad [ ] \quad MDP \quad Q \\
& \quad \quad \quad [ ] \quad MDP \quad MDP \\
& \quad \quad \quad MDP \quad MDP \\
& \quad \quad \quad [ ] \quad
\end{aligned}$$

---

<sup>1</sup> Multi-Agent Markov Decision Process

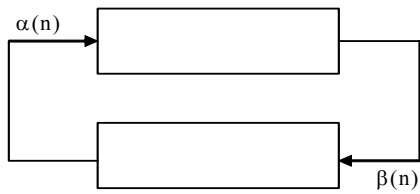
<sup>2</sup> Policy

<sup>3</sup> Online

<sup>4</sup> Greedy

<sup>5</sup> Interconnected Learning Automata

$$\begin{array}{llll}
& [\quad] . & & \\
\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\} & E \equiv \{\alpha, \beta, c\} & & \\
\beta & c \equiv \{c_1, c_2, \dots, c_r\} & \beta \equiv \{\beta_1, \beta_2, \dots, \beta_m\} & \\
& \beta_2 = 0 & \beta_1 = 1 & P \\
\beta(n) \ S & [\ , \ ] & \beta(n) \ Q & \\
c_i & \alpha_i & c_i \ . & [\ , \ ] \\
\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\} & \{\alpha, \beta, F, G, \phi\} & & \\
\phi(n) \equiv \{\phi_1, \phi_2, \dots, \phi_k\} & \beta \equiv \{\beta_1, \beta_2, \dots, \beta_r\} & & \\
G: \phi \rightarrow \alpha & F: \phi \times \beta \rightarrow \phi \ n & &
\end{array}$$



$$\begin{array}{lll}
\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\} & \{\alpha, \beta, p, T\} & \\
p = \{p_1, \dots, p_r\} & & \beta \equiv \{\beta_1, \beta_2, \dots, \beta_r\} \\
& & p(n+1) = T[\alpha(n), \beta(n), p(n)] \\
& & n \qquad \alpha_i
\end{array}$$

$$\begin{array}{l}
p_i(n+1) = p_i(n) + a[1 - p_i(n)] \\
p_j(n+1) = (1-a)p_j(n) \quad \forall j \neq i
\end{array} \tag{ }$$

$$\begin{array}{l}
p_i(n+1) = (1-b)p_i(n) \\
p_j(n+1) = (b/r-1) + (1-b)p_j(n) \quad \forall j \neq i
\end{array} \tag{ }$$

$$a : \quad b \quad a \quad b \quad a ( ) ( )$$

$$L_{RI} \quad b \quad L_{ReP} \quad a \quad b \quad L_{RP} \quad b$$

$$V(n) \quad (V(n)) \quad n$$

$$n \quad p(n) \quad (K(n))$$

$$p(n) \quad \alpha_i$$

$$p(n)$$

$$p_i(n) = prob[\alpha(n) = \alpha_i \mid V(n) \text{ is the set of active actions, } \alpha_i \in V(n)] = \frac{p_i(n)}{K(n)} \quad ( )$$

$$p_i(n+1) = p_i(n) + a.(1 - p_i(n)) \quad \alpha(n) = \alpha_i \quad ( )$$

$$p_i(n+1) = p_j(n) + a.p_i(n) \quad \alpha(n) = \alpha_i, \quad \forall j \quad j \neq i$$

$$p_i(n+1) = (1-b).p_i(n) \quad \alpha(n) = \alpha$$

$$p_i(n+1) = \frac{b}{r-1} + (1-b)p_j(n) \quad \alpha(n) = \alpha_i, \quad \forall j \quad j \neq i \quad ( )$$

$$: \quad p(n+1) \quad p(n)$$

$$p_j(n+1) = p_j(n+1).K(n) \quad \text{for all } j, \alpha_j \in V(n) \quad ( )$$

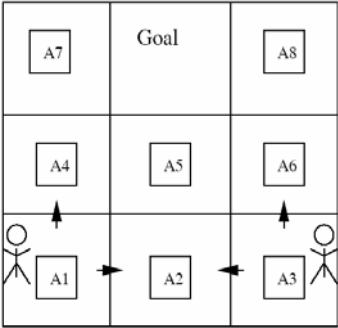
$$p_j(n+1) = p_j(n) \quad \text{for all } j, \alpha_j \notin V(n)$$

$$( )$$

$$( ) ( )$$

×

[ ] [ ]



(                    )

$$L_{RP}$$

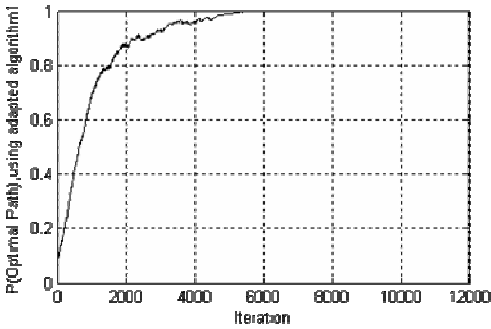
$j$

$$T_k^j$$

$$\begin{array}{c}
 \vdots \; k^j \\
 \\
 A_s \qquad \qquad s \qquad \qquad j \\
 \\
 a_m \qquad \qquad \qquad m \qquad \qquad (s,m) \\
 \\
 \qquad \qquad \qquad m \qquad \qquad j \\
 \\
 j \qquad \qquad \qquad s \leftarrow m \\
 \\
 L_{\pi_i}^j = R_G / t_{\pi_i}^j \qquad \qquad \pi_i \qquad \qquad j \\
 \\
 j \qquad \qquad \qquad t_{\pi_i} \qquad \qquad R_G \\
 \\
 \qquad \qquad \qquad \pi_i \\
 \\
 \qquad \qquad \qquad L_{\pi_i}^j < T_k^j \\
 \\
 -L_{\pi_i}^j \qquad \qquad L_{\pi_i}^j \\
 \\
 T_k^j \leftarrow T_{k-1}^j + \frac{1}{k^j} (L_{\pi_i}^j - T_{k-1}^j) \quad k^j \leftarrow k^j + 1 \\
 \\
 ( \qquad \qquad j \qquad \qquad )
 \end{array}$$

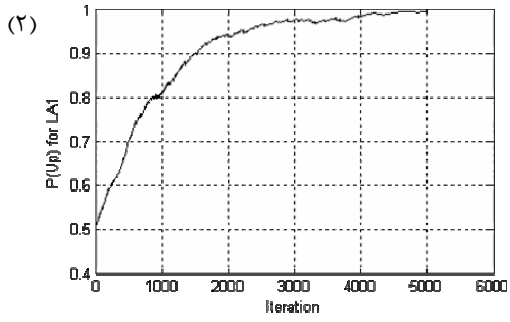
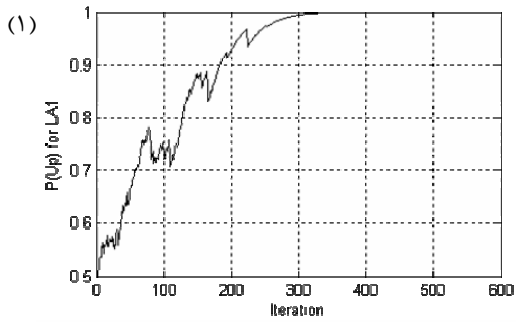
/

/



)

(



/ ( ) / ( )

:

$$L_{RP}$$

$j$

$$T_k^j$$

$$: k^j$$

$$A_s$$

$s$

$j$

$$a_m$$

$m$

$$(s,m)$$

$$A_s$$

( )

$$k^j \leftarrow k^j + 1$$

.( )

$m$

$j$

$$s \leftarrow m$$

$$L_{\pi_i}^j = R_G / t_{\pi_i}^j$$

$j$

$$\pi_i$$

$j$

$j$

$$t_{\pi_i}$$

$$R_G$$

$$\pi_i$$

$$L_{\pi_i}^j < T_k^j$$

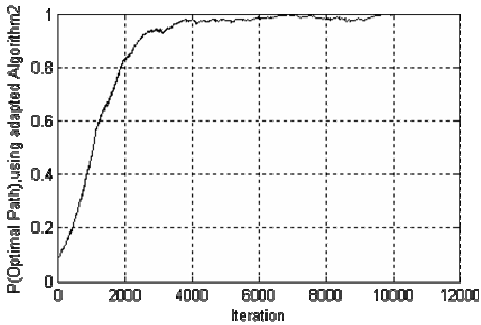
$$-L_{\pi_i}^j$$

$$L_{\pi_i}^j$$

$$T_k^j \leftarrow T_{k-1}^j + \frac{1}{k^j} (L_{\pi_i}^j - T_{k-1}^j) \quad k^j \leftarrow k^j + 1$$

$$(\hspace{1.5cm} j \hspace{1.5cm})$$

$$/ \hspace{1.5cm} /$$



$$[\hspace{0.5cm} ]$$

$$i \hspace{1.5cm} A_i$$

$$\begin{array}{l} j \hspace{1.5cm} k \hspace{1.5cm} r_j^i(k) \\ : \hspace{1.5cm} A_i \hspace{1.5cm} i \\ i \hspace{1.5cm} A_i \\ : \hspace{1.5cm} A_i \end{array}$$

$$\beta^i(n_i+1)=\frac{\rho_k^i(n_i+1)}{\eta_k^i(n_i+1)} \hspace{10cm} ( \hspace{0.5cm} )$$

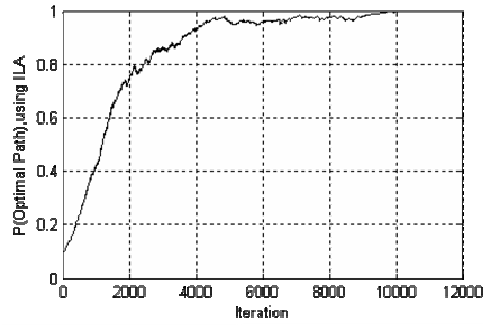
$$\eta_k^i(n_i+1) \hspace{0.5cm} i \hspace{1.5cm} k \hspace{1.5cm} \rho_k^i(n_i+1)$$

$$/$$

$$/$$

$$\wedge$$





	[ ]

/ ( / )

1. Weiss, G., "Multi-Agent Systems, a Modern Approach to Distributed Artificial Intelligence", MIT Press, 2000.
2. Nowe, A., Verbeeck, K., and Peeters, M., "Learning Automata as a Basis for Multi-Agent Reinforcement Learning", *First International Workshop on Learning and Adaptation in Multi-Agent Systems (LAMAS)*, pp. 71-85, 2006.
3. Hu, J. and Wellman, M. P., "Multi-Agent Reinforcement Learning: Theoretical Framework and an Algorithm", *Proceedings of the Fifteenth International Conference on Machine Learning*, 1998.
4. Laroche, P. Boniface, Y. and Schott, R., "A New Decomposition Technique for Solving Markov decision Processes", *Proceedings of the 2001 ACM Symposium on Applied Computing*, Las Vegas, Nevada, United States, pp. 12 – 16, 2001.
5. Peret, L. and Garcia, F., "On-line search for solving large Markov Decision Processes", *Sixth European Workshop on Reinforcement Learning (EWRL-6)*, 2003.
6. Nowe, A. and Verbeeck, K., "Colonies of Learning Automata", *IEEE Transactions on Systems, Man and Cybernetics*, vol. 32, pp. 772-780, 2002.

7. Beigy, H. and Meybodi, M. R., "Utilizing Distributed Learning Automata to Solve Stochastic Shortest Path Problem", *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 14, pp. 591-615, 2006.
8. Busoniu, L., De Schutter, B. and Babuska, R., "Learning and Coordination in Dynamic Multi-Agent Systems", *Delft Center for Systems and Control*, 2005.
9. Hanna, H., Mouaddid, A., "Markov Decision Process for Allocation Task in Multi-Agent Systems", *International Conference IPMU, France*, 2002.
10. Spaan, M., Vlassis, N. and Groen, F., "High Level Coordination of Agents based on Multi-Agent Markov Decision Processes with Roles", *IROS'02 Workshop on Cooperative Robotics*, 2002.
11. Shirazi, M. R. and Meybodi, M. R., "Application of Learning Automata to Cooperation in Multi-Agent Systems", *Proceedings of First International Conference on Information and Knowledge Technology (IKT2003)*, pp. 338-349, 2003.
12. Hetland, H. and Eriksen, T. L., "Interconnected Learning Automata Playing Iterated Prisoner's Dilemma", *Master Thesis, Agder University College*, 2004.
13. Verbeek, K., Nowe, A., Peeters, M. and Tuyls, K., "Multi-Agent Reinforcement Learning in stochastic Single and Multi-Stage Games", *Lecture Notes in Computer Science, Springer Berlin*, vol. 3394, pp. 275-294, 2005.