

# بهبود کارایی در سیستمهای چند عامله مبتنی بر اتوماتای یادگیر با استفاده از

## مفهوم آنتروپی

بهروز معصومی<sup>۱</sup>؛ محمد رضا میبدی<sup>۲</sup>

### چکیده

تا به حال برای مدل سازی سیستمهای چند عامله مدلهای مختلفی مبتنی بر مدل مارکوف پیشنهاد شده است که از جمله آنها مدل بازی های مارکوفی را می توان نام برد . در این مقاله روشی جدید با استفاده از مفهوم آنتروپی در سیستمهای چند عامله مبتنی بر اتوماتاهای یادگیر با هدف بهبود کارایی ارائه شده است. سیستم چند عامله مورد نظر، برای پیدا کردن خطمشی بهینه بازی های مارکوف مورد استفاده قرار می گیرد. در الگوریتم پیشنهادی، در هر حالت از محیط به ازای هر عامل بازی یک اتوماتای یادگیر قراردادده می شود تا بتوانند عاملهای محیط را کنترل نمایند. تعداد اعمال هر اتوماتای یادگیر با توجه به حالت های مجاور با آن تعیین می گردد و ترکیب اعمال اتوماتای یادگیر هر حالت از محیط، حالت بعدی را تعیین می کند. در الگوریتم مطرح شده اعمال انتخابی اتوماتاهای درون مسیر با توجه به هزینه مسیر طی شده و آنتروپی بدست آمده از بردارهای احتمالات اعمال اتوماتای یادگیر هر حالت، پاداش یا جریمه می گیرند. آزمایشهای انجام گرفته نشان داده اند که الگوریتم ارائه شده از کارایی مناسبی از نظر سرعت رسیدن به راهحل بهینه برخوردار است.

### کلمات کلیدی

سیستمهای چند عامله، اتوماتای یادگیر، بازی های مارکوف.

## Improving Learning Automata based Multi-agent System using Entropy Concept

Behrooz Masoumi<sup>1</sup>; Mohammad Reza Meybodi<sup>2</sup>

<sup>1</sup> Department of IT& Computer Engineering, Islamic Azad University, Qazvin Branch,  
Qazvin, Iran.

<sup>2</sup> Department of Computer Engineering and IT, Amirkabir University of Technology, Tehran, Iran  
mmeybodi@aut.ac.ir

### ABSTRACT

So far, many models have been proposed for multi-agent systems based on Markov Models. One of these models is Markov Games. In this paper a new method based on the entropy concept is presented to improvement of learning automata based multi-agent systems. Multi-agent system is used for finding optimal policies in Markov games.. In the proposed algorithm, each agent residing in every state of the environment is equipped with a learning automaton. Each joint-action of the learning automaton corresponds to moving to one of the adjacent states. Each agent moves from one state to another and tries to reach the goal state. The actions taken by the learning automata along the path traversed by the agent are then rewarded or penalized using the cost of the traversed path or the entropy of probability vector of learning automaton of each agent in the next state according to a learning algorithm. The results of experiments have shown that the proposed algorithms perform better than the existing algorithms in terms of cost and the speed of reaching the optimal policy.

### KEYWORDS

Markov Games, Learning Automata, Entropy.

<sup>1</sup> دانشکده مهندسی کامپیوتر و فناوری اطلاعات دانشگاه آزاد قزوین Email: masoumi@Qiau.ac.ir

<sup>2</sup> دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران. Email: mmeybodi@aut.ac.ir

یک سیستم چندعامله، در برگرنده جامعهای از عاملهای هوشمند و خود مختار است که در یک محیط در کنار یکدیگر سعی در انجام کاری خاص و رسیدن به هدفی مشخص دارند [۱]. امروزه در بسیاری از کاربردها و در زمینههای مختلف صنعتی، نظامی، مخابراتی، اطلاعاتی، از سیستمهای پیچیده و توزیعشده چندعامله استفاده فزاینده می شود [۲]. تاکنون برای مدل سازی سیستمهای چند عامله مدلهای مختلفی مبتنی بر مدل مارکوف پیشنهاد شده است که از جمله آنها مدل بازیهای مارکوفی<sup>۱</sup> را می توان نام برد [۳] [۴]. این بازی ها توسعه ای از فرآیندهای تصادفی مارکوف با چندین عامل بوده و به عنوان چارچوبی مناسب در تحقیقات یادگیری های چند عامله به ویژه یادگیری تقویتی چندعامله<sup>۲</sup> به کار رفته اند [۵] [۶]. بازیهای مارکوف، به صورت چند تایی (N,S,A,R,T) تعریف می گردند که در آن N تعداد عاملها، S مجموعه حالات، A مجموعه اعمال گروهی تمام عاملها، T تابع گذار حالات و R تابع پاداش با توجه به اعمال انتخابی عاملها در هر حالت تعریف می گردد. در این بازی ها، پاداش هر عامل وابسته به اعمال مشترک تمام عاملها بوده و حالت فعلی و گذار از حالتی به حالت دیگر از ویژگی مارکوف تبعیت می کنند. در حالت خاص، در صورتیکه که فقط یک عامل وجود داشته باشد بازی مارکوفی را فرآیند تصادفی مارکوف<sup>۳</sup> (MDP) گویند و در صورتیکه فقط یک حالت وجود داشته باشد به آن بازی های نرمال<sup>۴</sup> گفته می شود.

بازی های مارکوف دارای انواع مختلفی بوده و از نظر پاداش به دو رده بازی های کاملاً رقابتی (با مجموع صفر) و بازی جمع کلی تقسیم بندی شده اند. در بازی های رقابتی دوفره ساختار ماتریس پاداش هر عامل قرینه دیگری است. در بازی های مارکوف کلی، فرض می شود که برای پاداش عامل ها محدودیتی مطرح نیست. در حالت خاص، در صورتیکه پاداش یکسانی برای همه عاملها در نظر گرفته شود آنها را کاملاً همکارانه<sup>۵</sup> گویند و به آن فرآیندهای تصادفی مارکوف چند عامله (MMDP)<sup>۶</sup> نیز گفته می شود. در یک بازی های مارکوفی راه حل به معنای پیدا کردن سیاستی برای انتخاب اعمال توسط عاملها در هر حالت است تا بتواند امید ریاضی مجموع کاهشیافته پاداشها<sup>۷</sup> را برای همه عاملها بیشینه نماید. در MMDP ها با توجه به اینکه همه عاملها پاداش یکسان دریافت می کنند عامل ها بایستی یادگیرند تا در مورد سیاست بهینه توافق نمایند. در مقابل در بازی های مارکوفی کلی به دلیل وجود پاداش های متفاوت، پیدا کردن راه حل بهینه مشکل بوده و لذا نقاط تعادل<sup>۸</sup> (سیاست تعادل) در بازی مورد جستجو قرار می گیرند، وضعیتی که هیچ عاملی به تنهایی نمی تواند برای بهبود پاداشش سیاستش را تغییر دهد تا زمانی که تمام عاملهای دیگر سیاستشان را ثابت نگه می دارند.

برای پیدا کردن راه حل بهینه در بازی های مارکوفی، در شرایطی که شناخت کاملی از محیط وجود نداشته باشد و توابع پاداش و انتقال حالت مشخص نبوده و به صورت تجربی قابل مشاهده باشند استفاده از روشهای یادگیری تقویتی برای پیدا کردن نقاط تعادل مناسب است. در شکل کلی، این روش ها یک ورودی تابع انتخاب تعادل را دریافت نموده تا بتوانند با توجه شرایط تابع ارزش<sup>۹</sup> را محاسبه نموده و سیاست بهینه را یادگیرند. این روشها با توجه به پیدا کردن نقاط تعادل در هر مرحله از نظر محاسباتی پیچیدگی نسبتاً بالایی را دارند و در بسیاری از مواقع شرایط خاصی را بایستی در نظر گرفت. از جمله این روشها می توان به روش NASH-Q [۷]، PARETO-Q [۸] [۹] و NASH BARGAINING [۱۰] اشاره نمود.

اتوماتاهای یادگیر نیز در حال حاضر به عنوان ابزاری ارزشمند در طراحی الگوریتمهای یادگیری تقویتی بوده و واسطه ویژگیهایی که دارند در بسیاری از کاربردهای چند عامله و محیط های ناشناخته مناسب هستند [۱۱] [۱۲]. سادگی ساختار، نیاز به اطلاعات و بازخورد کم از محیط، قابلیت مناسب برای سیستمهای توزیع شده و سیستمهای چند عامله با اطلاعات ناکامل و ارتباطات محدود، انعطاف پذیری و قابلیت تحلیل در بیشتر کاربردها از جمله این ویژگیها هستند. برای حل بازی های مارکوفی نیز الگوریتمهای مختلفی مبتنی بر اتوماتاهای یادگیر ارائه شده است. در [۱۳] روشی مبتنی بر اتوماتاهای یادگیر برای حل فرآیندهای تصادفی مارکوف چند عامله و بازی های مارکوفی در شرایط ارگودیک<sup>۱۰</sup> مطرح شده اند و نشان داده شده است که شبکه ای از اتوماتاهای یادگیر<sup>۱۱</sup> قادر به رسیدن به استراتژی های تعادل در بازی های مارکوفی می باشند. در [۱۴] یک راه حل کلی برای بازی های مارکوفی با مجموع کلی با استفاده از اتوماتای یادگیر ارائه شده است که در آن با توجه به بردارهای احتمالات اعمال اتوماتای یادگیر و محاسبه آنتروپی اطلاعاتی به دست آمده در هر حالت، اعمال انتخابی اتوماتای یادگیر پاداش یا جریمه دریافت می دارند. در [۱۵] روشی برای حل بازی های مارکوفی همکارانه (MMDP) با استفاده از پیدا کردن هزینه کوتاهترین مسیر و پاداش دادن به اعمال اتوماتای مسیر بدست آمده ارائه شده است.

در این مقاله روشی مبتنی بر اتوماتاهای یادگیر برای حل بازی های مارکوف کلی برای خطمشی بهینه پیشنهاد شده است. در روش پیشنهادی، در هر حالت از محیط به ازای هر عامل، یک اتوماتای یادگیر قرار داده می شود تا بتواند عاملهای محیط را کنترل نمایند. با توجه به هزینه مسیر بدست آمده و نیز آنتروپی<sup>۱۲</sup> بردار احتمالات اتوماتای یادگیر حالت جدید اعمال انتخابی اتوماتاها پاداش می گیرند. برای بررسی و ارزیابی روش پیشنهادی از محیط بازی های Grid Game برای شبیه سازی آزمایش ها استفاده شده است. در ادامه مقاله، در بخش ۲ مفهوم بازی های مارکوفی توضیح داده شده و در بخش ۳ اتوماتاهای یادگیر و در بخش ۴ الگوریتم پیشنهادی مطرح و در بخش ۵ آزمایش های انجام شده در محیط بازی های Grid-Game و نتایج آزمایشها ارائه گردیده اند.

$$J(\alpha) = \lim_{l \rightarrow \infty} \frac{1}{l} E \left[ \sum_{t=0}^{l-1} R^{x(t)x(t+1)}(\alpha) \right] \quad (1)$$

## تعریف بازی های مارکوفی

بازی های مارکوف تعمیم فرآیندهای تصادفی مارکوف به حالت چندعامله است و بصورت زیر تعریف میشوند [۱۸]:  
**تعریف ۲.** بازی مارکوف بصورت چندتایی  $\langle n, S, A_1 \dots n, T, R \rangle$  بیان میشود که در آن  $n$  تعداد عامل ها،  $S$  مجموعه حالات،  $A_i$  مجموعه اعمال هر عامل  $i$  (در فضای اعمال گروهی  $A_1 \times A_2 \times \dots \times A_n$ )،  $T$  تابع انتقال  $[0, 1] \rightarrow S \times A \times S$  و  $R$  تابع پاداش برای عامل  $i$  نام با توجه به اعمال انتخابی در هر حالت است.  
 هر بازی مارکوفی با یک حالت بصورت یک بازی نرمال تکراری در تئوری بازی ها شناخته شده و هر بازی مارکوفی با یک عامل بصورت یک فرآیند تصمیم گیری مارکوفی است. علاوه بر این هر عامل تابع پاداش خاص خودش را داراست. در حالتی که هر عامل پاداش مختلفی را داراست پیدا کردن سیاست بهینه برای تمام عاملها بسیار مشکل بوده لذا بجای آن نقاط تعادل بازی جستجو می شوند، وضعیتی که هیچ عاملی به تنهایی نمی تواند تا زمانیکه تمام عاملهای دیگر سیاستشان را ثابت نگه می دارند برای بهبود پاداش سیاستش را تغییر دهد.

## ۲. اتوماتاهای یادگیر

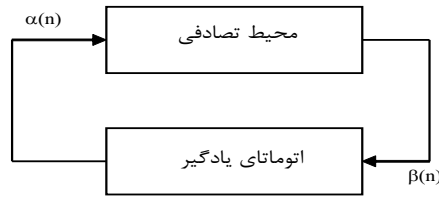
اتوماتای یادگیر، ماشینی است که میتواند تعدادی متناهی عمل را انجام دهد [۱۹]. هر عمل انتخاب شده توسط یک محیط احتمالی ارزیابی میشود و نتیجه ارزیابی در قالب سیگنالی مثبت یا منفی به اتوماتا داده میشود و اتوماتا از این پاسخ در انتخاب عمل بعدی تأثیر میگیرد. هدف نهایی این است که اتوماتا یاد بگیرد تا از بین اعمال خود، بهترین عمل را انتخاب کند. بهترین عمل، عملی است که احتمال دریافت پاداش از محیط را به حداکثر برساند. کارکرد اتوماتای یادگیر در تعامل با محیط، در شکل ۱ مشاهده می شود.  
 محیط را میتوان توسط سه تایی  $E \equiv \{\alpha, \beta, c\}$  نشان داد که در آن  $\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه ورودیها،  $\beta \equiv \{\beta_1, \beta_2, \dots, \beta_m\}$  مجموعه خروجیها و  $c \equiv \{c_1, c_2, \dots, c_r\}$  مجموعه احتمالات جرمیه می باشد. هرگاه  $\beta$  مجموعه های دو عضوی باشد، محیط از نوع  $P$  است. در چنین محیطی  $\beta_1 = 1$  به عنوان جرمیه و  $\beta_2 = 0$  به عنوان پاداش در نظر گرفته میشود. در محیط از نوع  $Q$ ،  $\beta(n)$  میتواند به طور گسسته یک مقدار از مقادیر محدود در فاصله  $[0, 1]$  را اختیار کند و در محیط از نوع  $S$ ،  $\beta(n)$  متغیر تصادفی در فاصله  $[0, 1]$  است.  $c_i$  احتمال اینکه عمل  $\alpha_i$  نتیجه نامطلوب داشته باشد. در محیط ایستا، مقادیر  $c_i$  بدون تغییر میمانند، حال آن که در محیط غیرایستا این مقادیر در طی زمان تغییر می کنند. اتوماتاهای یادگیر به دو دسته اتوماتای یادگیر با ساختار ثابت اتوماتای یادگیر با ساختار متغیر<sup>۳</sup> (VSLA) دسته بندی می شوند.  
 اتوماتای یادگیر با ساختار متغیر را میتوان توسط چهار تایی  $\{\alpha, \beta, p, T\}$  نشان داد که  $\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه عملهای اتوماتا،  $\beta \equiv \{\beta_1, \beta_2, \dots, \beta_r\}$  ورودیهای اتوماتا،  $p = \{p_1, \dots, p_r\}$  بردار احتمال انتخاب هریک از عملها و  $p(n+1) = T[\alpha(n), \beta(n), p(n)]$  الگوریتم یادگیری میباشد.  
 الگوریتم زیر براساس روابط (۲) و (۳) یک نمونه از الگوریتم های یادگیری خطی است. فرض میکنیم عمل  $\alpha_i$  در مرحله  $n$ ام انتخاب شود. در اینصورت پاسخ مطلوب از محیط بصورت:

$$\begin{aligned} p_i(n+1) &= p_i(n) + a[1 - p_i(n)] \\ p_j(n+1) &= (1-a)p_j(n) \quad \forall j \neq i \end{aligned} \quad (2)$$

و پاسخ نامطلوب از محیط بصورت زیر می باشد.

$$\begin{aligned} p_i(n+1) &= (1-b)p_i(n) \\ p_j(n+1) &= (b/r - 1) + (1-b)p_j(n) \quad \forall j \neq i \end{aligned} \quad (3)$$

در روابط (۲) و (۳)،  $a$  پارامتر پاداش و  $b$  پارامتر جرمیه میباشد. با توجه به مقادیر  $a$  و  $b$  سه حالت را می توان در نظر گرفت: اگر  $a$  و  $b$  باهم برابر باشند، الگوریتم را  $LRP$ ، هنگامیکه  $a$  از  $b$  خیلی کوچکتر باشد، الگوریتم را  $LR-EP$  و اگر  $b$  مساوی صفر باشد آن را  $LRI$  مینامیم [۱۹]. شمای  $S-LRP$  برای مدل های  $Q$  و  $S$  براساس رابطه (۴) بیان می شود:



شکل ۱- ارتباط بین اتوماتای یادگیر و محیط

اگر عمل  $\alpha_i$  در مرحله  $n$ ام انتخاب شود در این صورت طبق معادله (۴) داریم :

$$\begin{aligned} p_i(n+1) &= p_i(n) + a(1 - \beta_i(n))(1 - p_i(n)) \\ &\quad - a\beta_i(n)p_i(n) \\ p_j(n+1) &= p_j(n) + a(1 - \beta_i(n))(p_j(n)) + \\ &\quad a\beta_i(n)\left[\frac{1}{r-1} - p_j(n)\right] - a(1 - \beta_i(n))p_j(n) \text{ if } j \neq i \end{aligned} \quad (4)$$

$r$  تعداد اعمال ممکن،  $a$  پارامتر پاداش و  $b$  پارامتر جریمه میباشند. برای اطلاعات بیشتر در باره اتوماتاهای یادگیر می توان به [۲۰] مراجعه نمود.

## ۳. روش پیشنهادی

در این بخش روش پیشنهادی برای حل بازی های مارکوف کلی به کمک اتوماتای یادگیر بررسی می گردد. مسئله کنترل فرآیند تصمیمگیری مارکف (زنجیره مارکف) تکامل را میتوان با استفاده از شبکه های یادگیر متصل به هم مدل نمود که در آن، کنترل از یک اتوماتای یادگیر به اتوماتای یادگیر دیگر منتقل میشود [۲۱]. هر وضعیت از زنجیره مارکف دارای یک اتوماتای یادگیر است که طی مراحل فرآیند و با استفاده از روابط (۲) و (۳) سعی میکند توزیع احتمال بهینه اعمال در آن وضعیت را یاد بگیرد. عامل با شروع از وضعیت اولیه تا رسیدن به هدف بر روی این شبکه از اتوماتاهای یادگیر حرکت میکند و در هر وضعیت، از اتوماتای یادگیر در آن وضعیت برای انتقال به یکی از وضعیتهای مجاور کمک میگیرد که این کار با استفاده از برداراحتمال اعمال اتوماتای یادگیر انجام میگردد. در هر لحظه فقط یک اتوماتای یادگیر فعال بوده و انتقال از یک وضعیت به وضعیت دیگر، اتوماتای مربوط به وضعیت جدید را فعال مینماید. این فرآیند تا زمانی که بردار احتمالات کلیه اتوماتاهای یادگیر به پایداری برسد و یا شرط خاصی برقرار گردد تکرار میشود.

در یک بازی مارکوفی تغییر حالت با توجه به ترکیب اعمال مستقل انجام شده توسط اتوماتای یادگیر در هر حالت از محیط است. شبکه اتوماتای یادگیر به کار رفته برای MDP ها می تواند برای بازی های مارکوفی نیز با قرار دادن یک اتوماتای بازی هر عامل در هر حالت توسعه یابد [۱۳]. در الگوریتم پیشنهادی در هر حالت  $S_i$  ( $i=1..m$  و  $m$  تعداد حالات) از محیط بازی هر عامل  $k$  یک اتوماتای یادگیر نظیر  $LA(i, k)$  با ساختار متغیر قرار داده می شود. با توجه به تعداد حالات های مجاور با هر حالت از محیط، تعداد اعمال اتوماتای یادگیر مشخص می گردند. هر عامل با توجه به عمل انتخاب شده توسط اتوماتای یادگیر آن حالت و ترکیب گروهی اعمال انتخابی به حالت بعدی می رود. در ابتدا فرض می شود اتوماتاهای یادگیر تمام عملهای خود را با احتمالی یکسان انتخاب می کنند. در صورتیکه تغییر حالت ناشی از اعمال اتوماتای هر حالت منجر به ورود عامل به حالت نهایی (هدف) شود هزینه مسیر طی شده با هزینه مسیر قبلی مقایسه می شود اگر کمتر بود، هریک از اعمال انتخاب شده توسط اتوماتاهای یادگیر آن مسیر با توجه به میزان هزینه پاداش می گیرند و در غیر این صورت همه اتوماتاهای مسیر جریمه میشوند.

هزینه مسیر طی شده  $\pi_i$ ، طبق رابطه (۵) محاسبه می شود. در این رابطه،  $R_G$  پاداش مربوط به رسیدن به وضعیت هدف و  $t_\pi(k)$  مدت زمان سپری شده از شروع حرکت عامل  $k$  تا رسیدن آن به وضعیت هدف، با طی مسیر  $\pi$  می باشد.

در صورتیکه حالت جدید حالت نهایی (هدف) نیست هر اتوماتا با توجه به بردارهای احتمال اعمال اتوماتای یادگیر و محاسبه آنتروپی اطلاعاتی بردارهای احتمالات بصورت زیر تعریف می شود پاداش یا جریمه می گیرند.

$$L_k = \frac{t_\pi}{R_G} \quad (5)$$

$$I(S_i) = -\sum_{j=1}^N P_j(S_i) \log(P_j(S_i)) \quad (6)$$

$N$  تعداد اعمال اتوماتا در حالت  $S_i$  و  $P_i(s_i)$  احتمال انجام عمل  $j$  در حالت  $S_i$  است. آنتروپی بردار احتمال، میزان عدم قطعیت اتوماتای یادگیر حالت بعد را در انتخاب عمل خود نشان می‌دهد. هر چه آنتروپی بیشتر باشد میزان عدم قطعیت بیشتر است. در ابتدای الگوریتم که بردارهای احتمال اعمال اتوماتای یادگیر دارای مقادیر یکسان هستند یعنی  $P_1 = P_2 = \dots = P_r = 1/r$  میزان آنتروپی بیشترین مقدار است. لذا اتوماتای یادگیر دارای اطلاعات مفیدی برای رسیدن به هدف بوده و عملهای خود را به صورت تصادفی انتخاب می‌کند (جستجو ۱۴). با ادامه الگوریتم با تغییر احتمالات اعمال اتوماتای یادگیر میزان آنتروپی کم می‌شود و در حالت مینیم به صفر می‌رسد یعنی:  $\exists i, P_i = 1 \wedge \forall j \neq i P_j = 0$ . کاهش آنتروپی به این معنی است که اتوماتای یادگیر با احتمال بالایی یکی از اعمال خود را انتخاب می‌کند و دارای اطلاعات مفیدی برای رسیدن به هدف بوده و از این اطلاعات بهره برداری می‌نماید. برای اینکه مقدار آنتروپی به مقداری بین ۰ و ۱ تبدیل شده تا به عنوان پاداش یعنی  $\beta$  در اتوماتای یادگیر در محیطهای مدل  $S$  قابل استفاده باشد، آنتروپی آن حالت را به آنتروپی ماکزیمم تقسیم کرده و به آن آنتروپی نسبی گویند. الگوریتم نهایی در شکل ۲ نشان داده شده است.

## ۴. شبیه سازی انجام گرفته در محیط بازی Game-Grid

یکی از انواع بازیهای باتفاقی غیر رقابتی بازی های  $Game-Grid$  است که توسط Hu, Wellman ارائه شده است [۷]. این بازی یک بازی مارکوف دو نفری از نوع جمع کلی است. در شکل ۳ دو نوع بازی  $Grid Game$  نشان داده شده اند. در بازی نوع اول ( $GG1$ ) دو عامل وجود دارند که هردو می‌خواهند به یک هدف مشترک برسند و در بازی نوع دوم ( $GG2$ ) که یک بازی هماهنگی چند حالتی است، دو عامل و دو هدف وجود دارند. فرض بر این است که دو عامل از دو گوشه یک صفحه شروع کرده و سعی دارند تا با کمترین تعداد حرکت به هدف برسند. اعمال بازیکنان بطور همزمان انجام گرفته و هر یک از بازیکنان می‌توانند یکی از اعمال شمال، جنوب، شرق و یا غرب را انتخاب نمایند. مجموعه فضای حالات بصورت  $S = \{(0, 1), (0, 2), \dots, (8, 7)\}$  تعریف می‌شوند که هر حالت  $S = (l_r, l_p)$  مختصات عاملهای ۱ و ۲ را نشان می‌دهد. عاملها همزمان نمی‌توانند در یک مختصات یکسان قرار گیرند. اگر دو عامل سعی در حرکت به یک مربع یکسان داشته باشند حرکت هردو با شکست مواجه شده و هردو یک واحد جریمه گردیده و در موقعیت قبلی باقی می‌مانند. اگر عاملها به دو مربع مختلف غیر هدف بروند هر دو پاداش صفر را دریافت می‌کنند و اگر یکی به هدف برسد ۱۰۰ واحد پاداش می‌گیرد.

### The Proposed Algorithm

**Inputs:**  $a, b$ : reward and penalty parameter for each LA  $M$ : total training time

**Initialize :** In each state, a Learning Automaton of type S for each agent is placed. The set of actions of this LA is the set of permissible movements to other states.

**for all** states  $s$ , agents  $k$  **do**

$P(s,k) = 1/\text{number of permissible actions}$

$\text{TotalReward}(k) = 0$  // total reward received by each agent

$\text{LastAction}(s,k) = 0$  // last action played in each state

$T[k] = 0$ ; // Threshold for agent  $K$  ;

**end for**

//Main Loop

**for** episode = 1 to  $M$  **do**

$S = \text{StartState}$ ; //random or fixed

Iteration = 1

**while** not done

$\text{JointAction} = \emptyset$ ;

**for all** agent  $k$  **do** concurrently

Activate  $LA(S,k)$

Action = SelectAction ( $p(S, \text{agent})$ )

$\text{JointAction} = \text{JointAction} \cup \text{Action}$ ;

$\text{NewState} = \text{GetNextState}(S, \text{JointAction})$

$\text{LastAction}(\text{state}, \text{agent}) = \text{action}$ ;

$\text{TotalReward}(\text{agent}) = \text{TotalReward}(\text{agent}) + \text{GetReward}(S, \text{JointAction})$

StorePath( $\text{NewState}$ );

**if**  $\text{NewState} = \text{Goal state}$  **then**

**if**  $\text{CostLk} < TK$  **then**

reward  $LA(S,k)$  in Path with signal  $L_k$  based on Eq(5)

**else** penalize  $LA_k^s$  with  $(1.0 - L_k)$



```

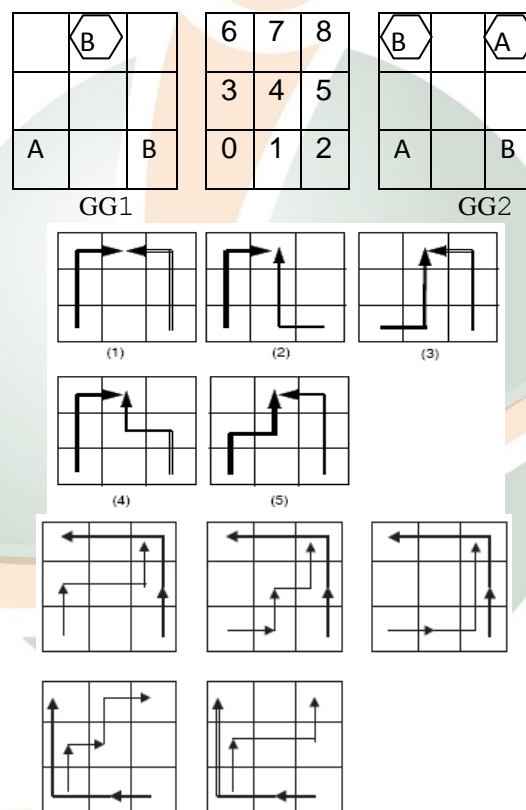

$$T_k \leftarrow T_k + \frac{1}{\text{Iteration}} (L_k - T_k)$$

endif
else Train  $LA(S, k)$  with reward signal Eq(6)
endif
end for
S=s',
Increment (Iteration)
end while
end for episode

```

شکل ۲. الگوریتم پیشنهادی

در هر دو بازی گذار از حالتی به حالت دیگر با قطعیت انجام می شود، یعنی حالت جاری و عمل مشترک عاملها منحصرأ حالت بعدی را تعیین می کنند.



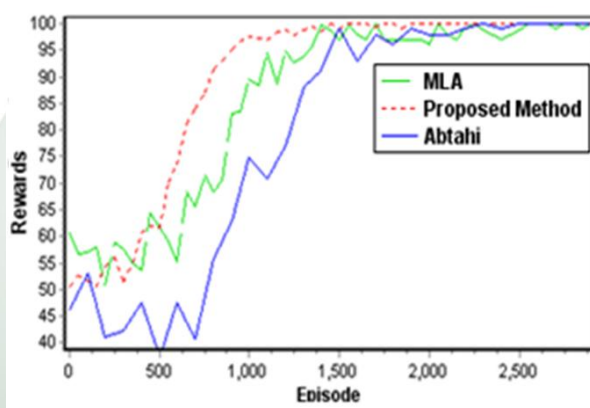
شکل ۳. طوی Grid World و نمایش مختصات بازی به همراه راه حل‌های بهینه در دوبازی

یک مسیر دنباله ای از اعمال انتخاب شده از نقطه شروع تا پایان را نشان می دهد. به کوتاهترین مسیری که یک عامل بدون تداخل با عاملهای دیگری می نماید مسیر بهینه گفته می شود. در این بازی فرض می شود عاملها از موقعیت هدف در ابتدای بازی آگاهی نداشته و همچنین از پاداش سراسری یکدیگر اطلاع ندارند. عاملها اعمالشان را همزمان انتخاب نموده و فقط می توانند از اعمال قبلی عاملهای دیگر و حالت فعلی (موقعیت مشترک هر دو عامل) آگاهی داشته باشند. برای حل مساله با توجه به روش ارائه شده در هر حالت از محیط بازی هر عامل یک اتوماتای یادگیر در نظر گرفته شده است. هر عامل با توجه به عمل مشترک گروهی ناشی از اعمال انتخابی اتوماتای یادگیر هر حالت به حالت جدید می رود و مسیر انتخابی نگهداری می شود. در صورتیکه حالت جدید حالت نهایی باشد هزینه مسیر طی شده با هزینه مسیر قبل مقایسه و در صورت بهبود به تمام اعمال انتخابی اتوماتاهای یادگیر آن مسیر پاداش داده می شود. در این مساله با توجه به اینکه حالت تکراری باعث افزایش طول مسیر می شود در صورتی که حالت جدید قبلاً مشاهده شده باشد اعمال انتخابی اتوماتاهای موجود جریمه تعلق می گیرد و در غیر اینصورت از آنتروپی بردار احتمالات اعمال اتوماتای یادگیر حالت جدید به عنوان پاداش یا جریمه اتوماتا استفاده می شود.

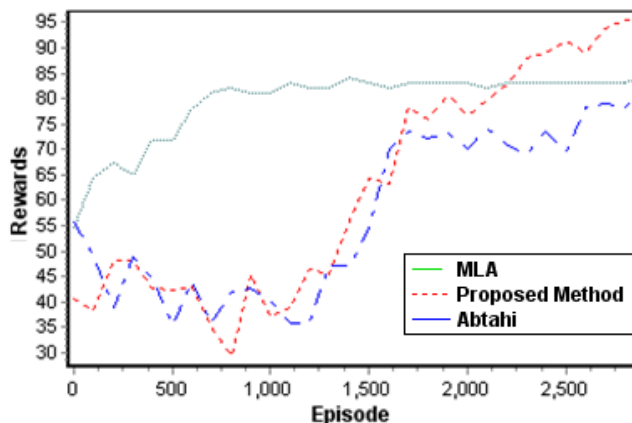
## مقایسه الگوریتم پیشنهادی با دیگر روشها

در ابتدا روش پیشنهادی در حل بازی های  $GG_1$  و  $GG_2$  با الگوریتمهای ارائه شده در [۱۴] یعنی الگوریتم MLA و [۱۵] Abtahi مقایسه می گردد. پارامترهای یادگیری را در تمام الگوریتمهای مطرح شده یکسان در نظر گرفته و از الگوریتم  $LRI$  در اتوماتای یادگیر استفاده می شود. آزمایشها به خاطر استحکام نتایج ۲۰۰ بار تکرار شده و در هر بار ۳۰۰۰ اپیزود آزمایش شده است. در ابتدا فرض می شود که نقطه شروع مختصات  $S=(0, 2)$  بوده و در هر اپیزود جدید عامل ها از همان نقطه مبدا مجددا شروع می کنند. برای بررسی میانگین پاداش تجمعی بدست آمده در هر اپیزود در ۲۰۰ آزمایش نشان داده شده است. شکل ۴ نمونه اجرای الگوریتمها را در بازی  $GG_1$  و شکل ۵ اجرای الگوریتمها در بازی  $GG_2$  نشان می دهد. همانطور که در شکل ۵ دیده می شود پس از تعدادی تکرار میانگین پاداش به سمت بیشترین افزایش می یابد. در آزمایش انجام گرفته پارامتر یادگیری  $a=0.01$  در نظر گرفته شده است و پاداش بدست آمده در هر اپیزود نشان داده شده است. همانطور که در شکل ۵ دیده می شود در هر دو بازی الگوریتم پیشنهادی رفتار بهتری را ارائه می دهد.

در آزمایش بعدی نقش پارامتر یادگیری در سرعت همگرایی مورد بررسی قرار می گیرد. همانطور که در شکلهای ۶ و ۷ دیده می شود، میزان پارامتر یادگیری نقش مهمی را در همگرایی الگوریتم ایفا می نماید. افزایش مقدار پارامتر  $a$  در بسیاری از مواقع باعث افزایش سرعت همگرایی می گردد. همچنین باید در نظر داشت در مواقعی که افزایش پارامتر یادگیری را داریم ممکن است باعث قرار گرفتن در نقاط بهینه محلی شده و باعث شود همگرایی بهینه دچار مشکل گردد. این مساله در بازی  $GG_2$  کاملا مشهود است. کاهش میزان پارامتر یادگیری اگرچه ممکن است نتایج بهتری ارائه دهد ولی باعث می شود تعداد تکرار ها برای رسیدن به همگرایی بیشتر گردد.

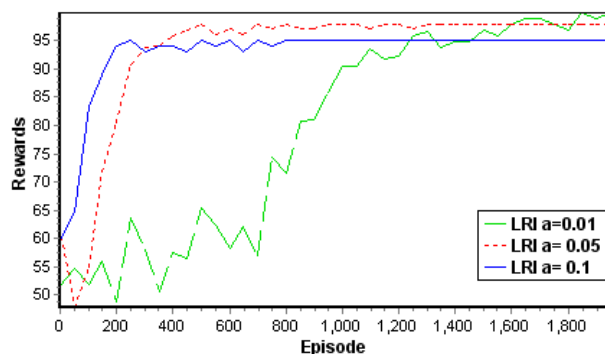


شکل ۴. مقایسه الگوریتم های مختلف بر روی مساله  $GG_1$  در ۳۰۰۰ اپیزود با پارامتر  $a=0.01, b=0$



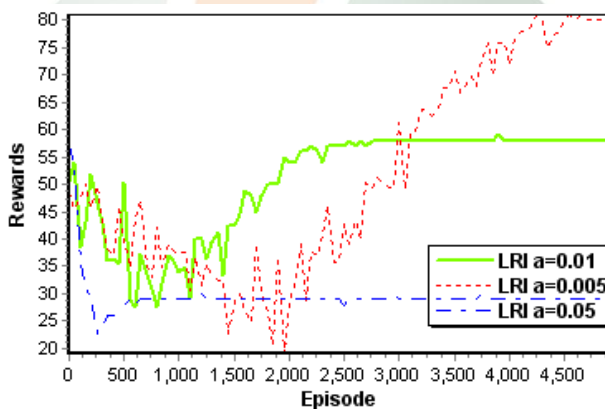
شکل ۵. مقایسه الگوریتم های مختلف بر روی مساله  $GG_2$

در ۳۰۰۰ اپیزود با پارامتر  $a=0.01, b=0$

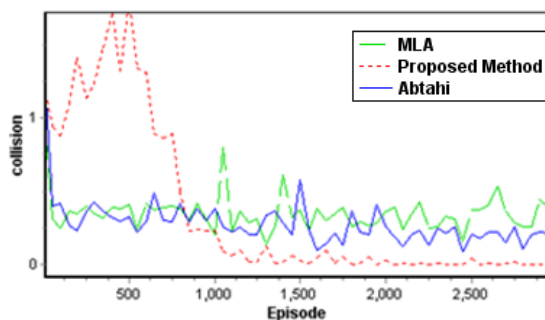


شکل ۶. بررسی رفتار پارامتر یادگیری در همگرایی الگوریتم در GG۱

برای بررسی رفتار الگوریتمها آزمایش دیگری انجام گرفت. که در آن متوسط تعداد برخوردهایی که بین دو عامل در هر اپیزود در ۲۰۰ آزمایش رخ می دهد محاسبه می شود. همانطور که در شکل ۸ دیده می شود در ابتدا الگوریتم پیشنهادی در حالت جستجو (Explore) بوده و به مرور طی فرآیند یادگیری به حالت استناد بر تجربه (Exploit) رسیده و لذا سریع همگرا می شود.



شکل ۷. بررسی رفتار پارامتر یادگیری در همگرایی الگوریتم در GG۲



شکل ۸. تعداد برخوردهای دو عامل در هر اپیزود برای بازی GG۱



## ۵. نتیجه گیری

در این مقاله روشی مبتنی بر اتوماتاهای یادگیر برای حل بازی های مارکوفی ارائه گردید. در این روش با توجه به هزینه مسیر طی شده تا رسیدن به هدف، اعمال انتخابی اتوماتاهای یادگیر در طول مسیر پاداش یا جریمه می گیرند. در طول مسیر نیز با توجه به آنتروپی به دست آمده از بردارهای احتمال اتوماتای یادگیر حالت جدید به عنوان پاداش های کمی جهت پاداش یا جریمه اتوماتا ها استفاده می شود و این کار باعث افزایش کارایی می گردد. با توجه به الگوریتم ارائه شده نتایج بدست آمده می بینیم روش مورد نظر رفتار مناسبتری را نسبت به روشهای دیگر برای حل بازی های مارکوفی نشان می دهد. تعداد تکرارها، پارامترهای یادگیری و سرعت رسیدن به نقاط تعادل را تعیین می نمایند. تنظیم پارامترهای پاداش و جریمه اتوماتاها می تواند کارایی رسیدن به راه حل بهینه را افزایش دهد. با توجه به نتایج به دست آمده اتوماتاهای یادگیر مدل مناسب یادگیری و هماهنگی بین عاملها در سیستمهای چندعامله بوده و می تواند به عنوان راه حلی مناسب و کارا در بازی های مارکوف به کار روند.

## ۶. مراجع

- [1] G. Weiss; Multiagent Systems: A Modern Approach to Distributed Artificial Intelligence, Cambridge, MA: MIT Press, 1999 .
- [2] A. K. Goel; V. Kummar and S. Irivasan; "Application of Multi-Agent System & Agent Coordination", 2nd National Conference Mathematical Techniques: Emerging Paradigms for Electronics and IT Industries(MATEIT-2008 ), New Dehli, India, 2008 .
- [3] Y. Shoham; Multiagent Systems: Algorithmic Game Theoretic and Logical Foundations, Cambridge University Press, 2009 .
- [4] M. L. Littman; "Markov Games as a Framework for Multi-agent Reinforcement Learning", In Proceedings of the 11th International Conference on Machine Learning, pp. 322 – 328, 1994 .
- [5] J. Osborne; A. Rubinstein. A Course in Game Theory, MIT Press, Cambridge, MA, 1994 .
- [6] L. Busni; R. Babuska and B. Schutter; "A Comprehensive Survey of Multiagent Reinforcement Learning ", IEEE Transaction on System, Man, Cybern, vol. 38, no.2, pp. 156–171, 2008 .
- [7] J. Hu and M. P. Wellman; "Nash Q-Learning for General-Sum Stochastic Games", Journal of Machine Learning Research, pp. 1039-1069, 2003 .
- [8] M. Song; G. Gu and G. Zhang; "Pareto-Q Learning Algorithm for Cooperative Agents in General-Sum Games", CEEMAS2005, pp. 576-578, ۲۰۰۵.
- [9] M. Song; J. Bai and R. Chen; "A New Learning Algorithm for Cooperative Agents in General-Sum Games", Proceedings of the Sixth International Conference on Machine Learning and Cybernetics, Hong Kong, 2007 .
- [10] H. Qio; F. Szidarovszky; Rozenblit and L. Yong; "Multi-agent Learning Model with Bargaining ", Proceedings of the 38th conference on winter simulation, pp. 934 - 940, 2006 .
- [11] M. R. Khojasteh and M. R. Meybodi; "Evaluating Learning Automata as a Model for Cooperation in Complex Multi-Agent Domains", Lecture Notes in Artificial Intelligence, Springer Verlag, LNAI 4434, pp. 409-416, 2007
- [12] A. Nowe; K. Verbeeck and M. Peeters; "Learning Automata as a basis for Multi-agent Reinforcement Learning", Lecture Notes in Computer Science, vol. 3898, pp. 71–85, 2006 .
- [13] P. Vrancx; K. Verbeeck and A. Nowe; "Decentralized Learning in Markov Games", IEEE Transactions on Systems, Man and Cybernetics (Part B: Cybernetics), vol. 38, iss. 4, pp. 976-81, 2008 .
- [14] B. Masoumi ; M. R. Meybodi and B. Jafarpour; "Solving General Sum Stochastic Games using Learning Automata", Proceedings of the second Joint Congress on Fuzzy and Intelligent Systems, Malek Ashtar University of Technology, Tehran, Iran, pp. 28-30, 2008 .
- [15] F. Abtahi and M. R. Meybodi; "Solving Multi-Agent Markov Decision Processes Using Learning Automata", Proceedings of the 6th International Symposium on Intelligent Systems (SISY2008 ), Subotica, Serbia, September 26-27, 2008 .
- [16] M. L. Puterman; Markov Decision Processes: Discrete Stochastic Dynamic Programming, John Wiley & Sons, Inc. New York, NY, USA, 1994 .
- [17] J. C. H. Watkins and P. Dayan; Q-Learning, Machine Learning. Vol. 3, pp.279–292, 1992 .
- [18] F. Thusijnsman; Optimality and Equilibria in Stochastic Games, (Centrum voor Wiskunde en Informatica, Amsterdam, 1992 ).
- [19] K. Narendra and M. Thathachar; Learning Automata: An Introduction, Prentice-Hall International, Inc, 1989 .
- [20] M. A. L. Thathachar, P. Sastry, "Varieties of Learning Automata: An Overview", IEEE Transaction on Systems, Man, and Cybernetics-Part B: Cybernetics, Vol. 32, No. 6, pp. 711-722, 2002 .
- [21] R. M. Wheeler and K. S. Narendra; "Decentralized Learning in Finite Markov Chains", IEEE Transactions on Automatic Control, pp. 519 – 526, 1986 .

- 
- ۱ Markov Game
  - ۲ Multi Agent Reinforcement Learning
  - ۳ Markov Decision Process
  - ۴ Normal Game
  - ۵ Fully Cooperative
  - ۶ Multi-Agent MDP
  - ۷ Sum of discounted expected rewards
  - ۸ Equilibrium Point
  - ۹ Value Function
  - ۱۰ Ergodic
  - ۱۱ Network of Learning Automata
  - ۱۲ Entropy
  - ۱۳ Variable Sturcture Learning Automata
  - ۱۴ Exploration
  - ۱۵ Exploitation



کنفرانس داده کاوی ایران