

# اتوماتاهای یادگیر، راه حلی برای بازی های غیر قطعی بامجموع کلی

بهروز معصومی<sup>۱</sup>، محمدرضا میبدی<sup>۲</sup>، برنا جعفرپور<sup>۲</sup>

<sup>۱</sup> دانشگاه آزاد اسلامی واحد علوم و تحقیقات، تهران و مرکز تحقیقات MRL دانشگاه آزاد اسلامی قزوین

<sup>۲</sup> دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران

Email: bmasoumi@Qazviniau.ac.ir, mmeybodi@aut.ac.ir, jafarpour@cic.aut.ac.ir

چکیده- بازی های غیر قطعی (اتفاقی) برای مدل سازی سیستمهای چند عامله بسیار مورد استفاده قرار گرفته اند. این بازیها توسعه ای از فرآیندهای تصادفی مارکوف با چندین عامل و بازی های ماتریسی با چندین حالت هستند. هدف هر عامل پیدا کردن سیاست بهینه ای است بطوریکه امید ریاضی مجموع کاهش یافته پاداشها را بیشینه نماید. در این مقاله یک مدل یادگیری تقویتی مبتنی بر اتوماتاهای یادگیر برای حل بازی های اتفاقی با مجموع کلی برای یافتن خط مشی بهینه پیشنهاد شده است. بازی هر حالت در محیط یک اتوماتا قرار داده شده بطوریکه تعداد اعمال هر اتوماتا با توجه به همسایگانش تعیین می گردد. هر اتوماتا مسوول انتخاب حالت بعدی محیط است. آزمایشهای انجام گرفته نشان داده اند که الگوریتم ارائه شده از کارایی مناسبی از هر دو جنبه هزینه و سرعت رسیدن به راه حل بهینه برخوردار است.

کلید واژه- سیستمهای چند عامله، اتوماتاهای یادگیر، یادگیری تقویتی چند عامله، بازی های اتفاقی

## ۱- مقدمه

بازی های غیر قطعی (اتفاقی)<sup>۵</sup> به عنوان توسعه ای از فرآیندهای تصادفی مارکوف با چندین عامل در سیستمهای چند عامله بسیار حائز اهمیت بوده و به عنوان چارچوبی مناسب در تحقیقات یادگیری های چند عامله به ویژه یادگیری تقویتی چندعامله (MARL)<sup>۶</sup> به کاررفته اند [۷،۶،۵]. یادگیری تقویتی چندعامله به سرعت در حال توسعه بوده و روشهای متنوع و مختلفی را درحوزه های رقابتی، همکاری و ترکیبی بر می گیرد بطوری که ارتباطی بین حوزه های مختلف علوم نظیر تئوری بازیها، بهینه سازی، یادگیری در بازی ها برقرار ساخته است [۸].

بازی های اتفاقی (SG) مدلی توسعه یافته از بازی های تکرارشونده هستند که در هر لحظه از زمان، بازی در یک حالت قرار دارد. گذار از حالتی به حالت جدید بر پایه تابع احتمالاتی با توجه به حالت قبلی و تعامل بین عامل ها در حالت قبل انجام می گیرد. هر حالت در یک SG می تواند بصورت یک فرآیند تصمیم گیری مارکوف دیده شود و هر SG با یک عامل بصورت یک فرآیند تصمیم گیری مارکوف می باشد. بازیهای اتفاقی دارای انواع متفاوتی

امروزه در بسیاری از کاربردها و در زمینه های مختلف صنعتی، نظامی، مخابراتی، اطلاعاتی، از سیستم های پیچیده و توزیع شده چندعامله استفاده فراوانی می شود [۲،۱]. یک سیستم چندعامله، در برگرنده جامعه ای از عامل های هوشمند و خود مختار است که در یک محیط درکنار یکدیگر در حال کار بوده و سعی در انجام کاری خاص و رسیدن به هدفی مشخص دارند. برای حل بسیاری از مسائل مهم دنیای واقعی مانند برخی از کاربردهای رباتیک، مسیریابی در شبکه، زمان بندی و تصمیم گیری اقتصادی نیازمند برنامه ریزی در حالت غیرقطعی هستیم. مدلهای فرآیند تصمیم گیری مارکوف، چارچوب مناسبی برای مدل سازی این مسائل و یافتن راه حل های بهینه برای آنها می باشد. برای مدل سازی سیستمهای چند عامله مدلهای مختلفی با توجه به مدل مارکوف پیشنهاد شده است که از جمله آنها مدل بازی های مارکوفی<sup>۲</sup>، فرآیندهای تصادفی مارکوف با مشاهدات جزئی<sup>۳</sup>، فرآیندهای تصادفی مارکوف با مشاهدات جزئی نامتمرکز<sup>۴</sup> را نام برد [۴،۳].

$$V(\pi, s) \equiv \sum_{t=0}^{\infty} \gamma^t E(r_t | \pi, s_0 = s) \quad (1)$$

این تابع، نگاشتی از مجموعه وضعیت‌ها به مقدار ارزش آن‌ها می‌باشد.  $r_t$  پاداش در زمان  $t$  و  $\gamma$  ضریب کاهش در محدوده  $[0, 1]$  است. بنابراین در اینجا هدف، یادگیری خط‌مشی بهینه است که به صورت زیر تعریف می‌شود:

$$\pi^*(s) = \arg \max_{\pi} [V(\pi(s, a))] \quad (2)$$

عامل باید سعی کند تا مقادیر  $V$  را به ازای خط‌مشی بهینه یاد بگیرد. این مقادیر را به طور خلاصه به شکل  $V^*$  نشان می‌دهیم.

$$\pi^*(s) = \arg \max_a [r(s, a) + \gamma V^*(\delta(s, a))] \quad (3)$$

برای محاسبه تابع  $V(s)$  از تکنیک برنامه‌سازی پویا<sup>۸</sup> استفاده می‌شود. در این صورت، روشی تقریبی برای تخمین مقادیر بهینه  $V(s)$  به کار می‌رود که روش *Value Iteration* نام دارد. این روش هنگامی قابل استفاده است که عامل، توابع  $\delta: S \times A \rightarrow S$  و  $r: S \times A \rightarrow R$  را بشناسد، در غیراین صورت نمی‌توان این روش را به کار گرفت. در چنین حالتی، از الگوریتمی به نام یادگیری  $Q$  استفاده می‌نماییم. در یادگیری  $Q$  استاندارد یک عامل مقادیر  $Q$  را طبق معادله زیر یاد می‌گیرد که در آن  $\alpha$  نرخ یادگیری و  $\beta$  ضریب کاهش است.

$$Q_{k+1}(s_k, a) = (1 - \alpha_t) Q_k(s_k, a) + \alpha_t \left[ r_t + \beta \max_b Q_k(s_{k+1}, b) \right] \quad (4)$$

بازی های اتفاقی تعمیم فرآیند تصادفی مارکوف به حالت چندعامله و همچنین توسعه‌ای از بازی‌های ماتریسی با چندین حالت بوده و به آنها نام بازیهای مارکوفی نیز گویند.

**تعریف ۲.** یک بازی اتفاقی (مارکوف) بصورت چندتایی  $\langle n, S, \vec{A}, \vec{R}, T \rangle$  بیان میشود که در آن  $S$  مجموعه حالات،  $n$  تعداد عامل ها،  $\vec{A}$  مجموعه اعمال هر عامل  $i$ ،  $\vec{R}$  تابع پاداش عامل های مختلف  $R: S \times \vec{A} \rightarrow R$  و  $T$  تابع تبدیل اتفاقی است. در حالت چند عامله تابع تبدیل با توجه به عمل گروه عاملها تعیین می شود. اگر پاداش تمام

هستند. بازی های اتفاقی از نظر پاداش به بازی های با مجموع صفر (رقابتی) و مجموع کلی تقسیم بندی شده اند. برای حل بازی های اتفاقی اعم از بازی های رقابتی و غیر رقابتی الگوریتمهای یادگیری تقویتی متعددی به کار رفته اند. با توجه به پیچیدگی محاسباتی روش ها هر کدام در کاربرد خاص با شرایط خاصی استفاده شده اند. هدف اصلی این مقاله ارائه الگوریتمی با استفاده اتوماتاهای یادگیر برای حل مسائل بازی های اتفاقی مجموع کلی و یافتن خط‌مشی بهینه است. در بخش ۲ به مبحث یادگیری تقویتی و بازی های اتفاقی پرداخته شده در بخش ۳ اتوماتاهای یادگیر و در بخش ۴ و ۵ الگوریتم پیشنهادی و نتایج آزمایشها ارائه گردیده اند.

## ۲- یادگیری تقویتی و حل بازی های اتفاقی

در سیستمهای چند عامله، الگوریتمهای یادگیری تقویتی مانند یادگیری  $Q$ ، با موفقیت در بسیاری از کاربردها مورد استفاده قرار گرفته‌اند [۸]. اکثر الگوریتمهای یادگیری تقویتی چندعامله، بر پایه روشهای تک‌عامله بنا نهاده شده‌اند. در حالت یادگیری تک عامله که در آن عامل به طور مستقل در حال یادگیری باشد، فرآیند یادگیری تقویتی را می‌توان به صورت فرآیند تصمیم‌گیری مارکوف تعریف نمود.

**تعریف ۱.** فرآیند تصادفی مارکوف بصورت چندتایی  $\langle S, A, P, r, \gamma \rangle$  نشان داده می شود که در آن  $S$  مجموعه متناهی از وضعیت‌ها؛  $A$  مجموعه عملیات قابل دسترس برای عامل،  $\gamma$  ضریب کاهش و  $P: S \times S \times A \mapsto [0, 1]$  احتمال انتقال از وضعیت جاری به وضعیت بعدی با انجام عمل  $a$  است و  $r: S \times A \mapsto \mathbb{R}$  تابع پاداش است که یک مقدار عددی را بر می گرداند.

در یک  $MDP$  هدف عامل پیدا کردن استراتژی  $\pi: S \rightarrow A$  است که امید ریاضی مجموع کاهش یافته پاداشها را بیشینه نماید. یادگیری تقویتی برپایه اصل بهینگی *Bellman* استوار است. برای هر خط‌مشی  $\pi$  که عامل می‌تواند دنبال کند، بر روی وضعیت‌ها تابعی به نام تابع ارزیابی<sup>۹</sup> به شکل زیر تعریف می‌شود:

استراتژی و اعمال دیگران دارد. بطورکلی معادله بروز رسانی  $Q$  را می توان بصورت زیر نوشت :

$$Q_i(s, \vec{a}) = (1 - \alpha) Q_i(s, \vec{a}) + \alpha (r + \gamma V(s')) \quad (8)$$

$$V(s) = Value([Q_i(s, \vec{a})])$$

در این معادله  $r_t^i$  پاداش عامل  $i$  در لحظه  $t$  و  $\gamma$  ضریب کاهش است. الگوریتم یادگیری  $Q$  چند عامله در شکل (۱) دیده می شود: انتخاب توابع مختلف متناسب با نوع بازی باعث ایجاد الگوریتمهای مختلف شده است. از جمله این توابع می توان به  $minimaxQ$  برای بازی های رقابتی [۵]،  $Nash(Q)$  [۷] و  $NashBargainingQ$  [۹] و  $ParetoQ$  [۱۰] برای بازیهای با مجموع کلی استفاده شده است .

Multi-Q-learning (MarkovGame,  $\alpha, \gamma$ )  
Inputs: discount factor  $\gamma$ , learning rate  $\alpha$   
func  $Q_i^*$  Output : state-value func  $V_i^*$ , action-value  
Initialize :  $s, a_1..a_n, Q_1..Q_n$   
1. for  $i=1$  to  $M$   
2. simulate actions  $a_1..a_n$  in state  $s$   
3. observe rewards  $R_1..R_n$  and Next State  $s'$   
4. for  $i = 1$  to  $N$   
(a) compute  $V_i(s')$   
5. (b)  
 $Q_i(s, \vec{a}) = (1 - \alpha) Q_i(s, \vec{a}) + \alpha (1 - \gamma) (R_i + \gamma V_i(s'))$   
agents choose action  $a'_1..a'_n$   
6.  $s = s', a_1 = a'_1, a_n = a'_n$   
7. decay  $\alpha$

شکل ۱: الگوریتم یادگیری  $Q$  چند عامله

### ۳- اتوماتاهای یادگیر

با توجه به اینکه مدل های اتوماتاهای یادگیر ارتباطی نزدیک با بحث یادگیری تقویتی چند عامله دارند، در این بخش اتوماتاهای یادگیر به عنوان مدلی از یادگیری تقویتی به اختصار شرح داده میشود.

#### ۳-۱- اتوماتاهای یادگیر

اتوماتای یادگیر، ماشینی است که می تواند تعدادی متناهی عمل را انجام دهد. هر عمل انتخاب شده توسط یک محیط احتمالی ارزیابی می شود و نتیجه ارزیابی در قالب سیگنالی مثبت یا منفی به اتوماتا داده می شود و اتوماتا از این پاسخ در انتخاب عمل بعدی تأثیر می گیرد.

عاملها برابر باشد بازی را کاملاً همکارانه<sup>۹</sup> و اگر  $n=2$  و پاداش یک بازیکن مخالف پاداش دیگری باشد بازی را رقابتی و در حالت کلی بازی را ترکیبی گویند.

در حالت یادگیری تقویتی چند عامله ، بیشینه نمودن سودمندی (پاداش) مورد انتظار هر عامل به تنهایی کافی نیست بطوریکه هدف پیدا کردن سیاست متعادل در بازی های مارکوف است. از جمله این سیاست ها می توان سیاست تعادل نش<sup>۱۰</sup> را نام برد. به عبارت دیگر پیدا کردن یک سیاست متعادل به عنوان یک راه حل برای بازی های اتفاقی محسوب می شود. استراتژی یک عامل، برنامه و طرح آن در بازی محسوب می گردد. استراتژی بصورت  $\pi = (\pi_0, \dots, \pi_t, \dots)$  روی کل بازی تعریف می شود که  $\pi_t$  قانون تصمیم گیری در لحظه  $t$  نامیده می شود. قانون  $\pi$  را ایستا گویند اگر مستقل از زمان باشد. در این مقاله ما روی استراتژی های ایستا تأکید داریم. استراتژی بصورت معادله ۵ تعریف می شود که توزیع احتمالاتی را بر روی اعمال قابل دسترس عاملها برای هر حالت  $m..1, j, s_j$  و  $m$  تعداد حالات است تعیین می کند.

$$\bar{\pi} = (\bar{\pi}(s_1), \dots, \bar{\pi}(s_m)) \quad (5)$$

برای هر حالت  $s_j$  طبق معادله ۶ داریم، بطوریکه  $P(a_k)$  احتمال انجام عمل  $a_k$  را نشان می دهد .

$$\bar{\pi}(s_j) = \{P(a_1), \dots, P(a_m)\} \quad (6)$$

$$\sum_{k=1}^M P(a_k) = 1 \text{ و } M \text{ تعداد اعمال هر عامل است}$$

تعریف ۳. در یک بازی اتفاقی یک تعادل نش بصورت زوج استراتژی  $(\pi_1^*, \pi_2^*)$  تعریف می شود بطوریکه برای هر  $s \in S$  داریم :

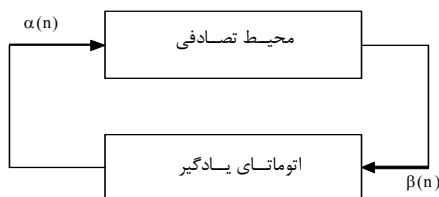
$$v^1(s, \pi_1^*, \pi_2^*) \geq v^1(s, \pi_1, \pi_2^*) \forall \pi_1 \in \Pi^1 \quad (7)$$

$$v^2(s, \pi_1^*, \pi_2^*) \geq v^2(s, \pi_1^*, \pi_2) \forall \pi_2 \in \Pi^2$$

یادگیری  $Q$  چند عامله مشابه یادگیری تک عامله است با این تفاوت که در اینجا اعمال و استراتژی های گروهی<sup>۱۱</sup> مد نظر است و پیدا کردن مقدار  $Q$  بهینه بستگی به

$$\begin{aligned}
 p_i(n+1) &= (1-b)p_i(n) \\
 p_j(n+1) &= (b/r-1) + (1-b)p_j(n) \\
 \forall j \quad j \neq i
 \end{aligned} \quad (10)$$

در روابط (۹) و (۱۰)، پارامتر پاداش و  $b$  پارامتر جریمه می‌باشند. با توجه به مقادیر  $a$  و  $b$  سه حالت را می‌توان در نظر گرفت: اگر  $a$  و  $b$  با هم برابر باشند، الگوریتم را  $L_{RP}$ ، هنگامی که  $b$  از  $a$  خیلی کوچکتر باشد، الگوریتم را  $L_{REP}$  و اگر  $b$  مساوی صفر باشد آن را  $L_{RI}$  می‌نامیم [11]. شمای  $SLR-P$  برای مدل‌های  $Q$  و  $S$  براساس رابطه (۱۱) بیان می‌شود:



شکل ۲- ارتباط بین اتوماتای یادگیر و محیط

اگر عمل  $\alpha_i$  در مرحله  $n$ ام انتخاب شود در این صورت طبق معادله (۱۱) داریم:

$r$  تعداد اعمال ممکن، پارامتر پاداش و  $b$  پارامتر جریمه

$$\begin{aligned}
 p_i(n+1) &= p_i(n) + a(1-\beta_i(n))(1-p_i(n)) \\
 &\quad - a\beta_i(n)p_i(n) \\
 p_j(n+1) &= p_j(n) - a(1-\beta_i(n))p_j(n) + \\
 &\quad a\beta_i(n)\left[\frac{1}{r-1} - p_j(n)\right] \quad j \neq i
 \end{aligned} \quad (11)$$

می‌باشند. برای اطلاعات بیشتر در باره اتوماتاهای یادگیر می‌توان به [12] مراجعه نمود.

#### ۴- استفاده از اتوماتاهای یادگیر برای حل بازی های اتفاقی

در این بخش روش پیشنهادی برای یادگیری در سیستمهای چند عاملی به کمک اتوماتای یادگیر را بررسی می‌کنیم. ما این روش را  $MLA$ <sup>۱۳</sup> می‌نامیم. در یک بازی مارکوفی عمل انتخاب شده در هر حالت با توجه به نتیجه اعمال گروهی عاملهای مستقل در سیستم است. در این مدل در هر حالت  $s_i$  ( $m \leq i \leq m$ ) تعداد حالات از محیط بازی هر عامل  $k$  یک اتوماتای یادگیر نظیر  $LA_k^i$

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 ۲ مشاهده می‌شود.

محیط را می‌توان توسط سه تایی  $E \equiv \{\alpha, \beta, c\}$  نشان داد که در آن مجموعه  $\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه ورودی‌ها،  $\beta \equiv \{\beta_1, \beta_2, \dots, \beta_m\}$  مجموعه خروجی‌ها و  $c \equiv \{c_1, c_2, \dots, c_r\}$  مجموعه احتمال‌های جریمه می‌باشد. هرگاه  $\beta$  مجموعه‌ای دو عضوی باشد، محیط از نوع  $P$  است. در چنین محیطی  $\beta_1 = 1$  به عنوان جریمه 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 می‌باشد. در محیط ایستا، مقادیر  $c_i$  بدون تغییر می‌مانند، حال آن‌که در محیط غیرایستا این مقادیر در طی زمان تغییر می‌کنند. اتوماتاهای یادگیر به دو دسته اتوماتای یادگیر با ساختار ثابت اتوماتای یادگیر با ساختار متغیر ( $VS LA$ )<sup>۱۲</sup> دسته بندی می‌شوند.

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 (۹) و (۱۰) یک نمونه از الگوریتم‌های یادگیری خطی است. فرض می‌کنیم عمل  $\alpha_i$  در مرحله  $n$ ام انتخاب شود.

- پاسخ مطلوب از محیط

$$\begin{aligned}
 p_i(n+1) &= p_i(n) + a[1 - p_i(n)] \\
 p_j(n+1) &= (1-a)p_j(n) \quad \forall j \quad j \neq i
 \end{aligned} \quad (9)$$

- پاسخ نامطلوب از محیط

با ساختار متغیر و مدل  $S$  قرار داده می‌شود. با توجه به تعداد همسایگان هر اتوماتا، تعداد اعمال اتوماتا در هر حالت در نظر گرفته شده و هر عمل مشترک متناظر با انتقال به یکی از حالات همسایه می‌باشد. هر عامل برای تعیین حالت بعدی خود از اتوماتای یادگیر با توجه به عمل انتخابی عاملهای محیط کمک می‌گیرد.

در ابتدا اتوماتاهای یادگیر تمام عملهای خود را با احتمالی یکسان انتخاب می‌کنند. در صورتیکه عمل اتوماتا هامنجر به ورود عامل به حالت هدف شود اتوماتا پاداش می‌گیرد، در غیر اینصورت از آنتروپی بردار احتمال اتوماتای یادگیر حالت بعد برای تعیین پاداش یا جریمه استفاده می‌شود. 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. عدم قطعیت بالا در بردار احتمال اتوماتای یادگیر به این معنی است که این اتوماتا دارای اطلاعات مفیدی برای رسیدن به هدف نیست و عملهای خود را به صورت تصادفی انتخاب می‌کند (جستجو<sup>۱۴</sup>). ولی چنانچه عدم قطعیت کم باشد به این معنی است که اتوماتا با احتمال بالایی یکی از اعمال خود را انتخاب می‌کند و دارای اطلاعات مفیدی برای رسیدن به هدف می‌باشد و از این اطلاعات بهره برداری می‌نماید<sup>۱۵</sup>. فرض کنید که  $\{p_1, p_2, \dots, p_r\}$  بردار احتمال اعمال یک اتوماتای یادگیر باشد. آنتروپی این بردار احتمال به شکل زیر تعیین می‌شود.

$$E(P) = -\sum_{i=1}^r p_i \log(p_i) \quad (12)$$

زمانی آنتروپی بیشترین مقدار را خواهد داشت که تمام اعمال احتمالی یکسان داشته باشند  $P_{equal} = \{p_1 = p_2 = \dots = p_r = 1/r\}$  و زمانی کمترین مقدار (برابری) را خواهد داشت که  $\exists i, p_i = 1 \wedge \forall j \neq i, p_j = 0$  برای اینکه مقدار آنتروپی را به مقداری بین ۰ و ۱ تبدیل کنیم تا به عنوان بردار تقویتی در اتوماتای ساختار متغیر مدل  $S$  قابل استفاده باشد از فرمول ۱۳ استفاده می‌شود [۱۳].

$$\beta = E((P) / E(P_{equal}))^K \quad (13)$$

$K$  پارامتر روش می‌باشد. میزان این پارامتر باید به دقت تعیین شود. مقادیر بالای این پارامتر باعث می‌شود که  $\beta$  ها مقادیر کمی داشته باشند و اتوماتاهای یادگیر بیش از حد پاداش ببینند. این به معنی جستجوی بیشتر در محیط است. هر چه میزان  $K$  کمتر باشد  $\beta$  ها مقادیر بیشتری خواهند داشت. این امر باعث می‌شود که اتوماتاها بیشتر جریمه شوند. این مسئله باعث می‌شود که حتی حالت‌های مطلوب نیز پاداش لازم را نگیرند. در الگوریتم پیشنهادی به گونه ای عمل می‌شود مشابه با روش بولتزمن در یادگیری  $Q$  در اینجا ابتدا عاملها جستجو انجام داده و به مرور زمان با توجه به تغییر پارامتر بهره برداری از دانسته های خود داشته باشند. فرض کنید که عامل  $k$  در حالت  $s$  باشد و اتوماتای یادگیر آن  $LA_k^s$  عامل را به حالت  $s'$  هدایت کند. در اینصورت تعیین سیگنال تقویتی طبق معادله ۱۴ تعیین می‌شود بطوریکه  $P(LA_i^s)$  نشان دهنده بردار احتمال اتوماتای یادگیر  $i$  قرار گرفته در حالت  $s$  می‌باشد. الگوریتم  $MLA$  در شکل (۳) آورده شده است.

#### **MAL (StochasticGame,a,b,k,M)**

Inputs:  $a, b$ : reward and penalty parameter for LA  
 $k$ : exploration Parameter ,  $M$ : total training time  
 Initialize :  $s_0, a_1, \dots, a_n$ , initialization probability of all of LA

1. For Episod = 1 to M do
2. While not done
3.  $k = \text{initial Value}$
4. for each agent  $k$  do concurrently
5. Active  $LA_k^i$
6. Choose action  $a_k^i$  in state  $s_i$
7. Observe Rewards  $r_k^i$  and Next State  $s'$
8. Compute  $\beta_k^i$  signal based on EQ (14)
9. Train  $LA_k^i$  residing in state  $S_i$  according  $\beta_k^i$
10.  $s = s'$
11. End while
12. Increment  $k$
13. End for

شکل ۳- الگوریتم MLA

$$\beta_i^s = \begin{cases} 1 & s' = \text{Out of bound} \\ & s' = \text{Goal} \\ (E(P(LA_i^{s'})) / E(P_{equal}))^K & \text{otherwise} \end{cases} \quad (14)$$

#### ۴-۱- حل بازی های اتفاقی با اتوماتای یادگیر

یکی از انواع بازی های اتفاقی غیر رقابتی بازی های *Grid Game* است که توسط *Hu, Wellman* ارائه شده است [۶].

این بازی یک بازی دو نفری از نوع جمع کلی است. در این بازی دو عامل از دو گوشه یک صفحه شروع کرده و سعی می کنند تا به مربع هدف برسند. این بازی در چند شکل مختلف مطرح شده است. در یک نوع آن فقط یک هدف وجود دارد و و دارای ۹ خانه است. حرکت بازیکنان قطعی بوده و هدف این است که با کمترین تعداد حرکت به هدف برسیم. شکل (۴) مدل بازی رابه همراه مختصات مورد نظر بازی و همچنین راه حل های بهینه را نشان می دهد.

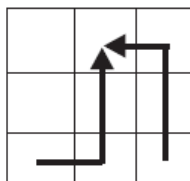
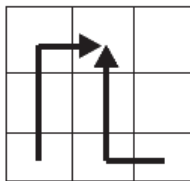
یک مسیر دنباله ای از اعمال از نقطه شروع تا پایان را نشان می دهد. در اصطلاح بازی چنین مسیری را استراتژی یا سیاست می نامند. کوتاهترین مسیری که با مسیر عامل دیگر تداخل نداشته باشد را استراتژی بهینه می نامند بطوریکه یک موازنه نش را می سازند زیرا هر مسیر (استراتژی) بهترین پاسخ در قبال دیگری است. در این بازی فرض می شود عاملها از موقعیت هدف در ابتدای بازی آگاهی نداشته و همچنین از پاداش کلی 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.

حالت دیگر با قطعیت انجام می شود. یعنی حالت جاری و عمل مشترک عاملها منحصرأ حالت بعدی را تعیین می کنند.

برای حل مساله با توجه به روش ارائه شده بازی هر حالت برای هر عامل یک اتوماتای یادگیر در نظر گرفته شده است. هر عامل با توجه به عمل مشترک گروهی از یک حالت با حالت جدید می رود (در صورت عدم برخورد دو عامل با یکدیگر) با توجه به عدم برخورد با موانع (دیوارها) و براساس معادله ۱۴ مقدار  $\beta$  (پارامتر پاداش یا جریمه) تعیین شده و اتوماتای آن حالت بروز می گردد. عاملها اعمالشان را همزمان انجام داده و هر دو قابلیت مشاهده حالت جدید، پاداش آنی بدست آمده و عمل انجام شده توسط دیگری را دارند.

	Goal	
A1		A2

6	7	8
3	4	5
0	1	2



شکل ۴. بازی Grid World و نمایش مختصات بازی به همراه راه حل های بهینه

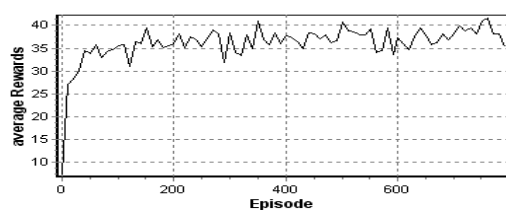
#### ۵- آزمایشهای انجام گرفته و نتایج آنها

هدف اصلی آزمایشها بررسی این است که آیا اتوماتای یادگیر می تواند برای راه حل بازی های مارکوفی مناسب باشد؟ آیا به راه حل تعادل نش بهینه می شود؟ در آزمایش اولیه بررسی همگرایی اتوماتا انجام گرفته است. تغییرات احتمال انتخاب عمل حرکت به بالا در خانه شروع با توجه به پارامترهای در نظر گرفته شده یعنی پارامتر  $a$  برابر  $0.02$  و پارامتر  $b$  برابر  $0.002$  در عامل ۱ بررسی شد. شکل ۵ نمودار تغییرات را برای این عامل نشان می دهد. با توجه به پارامترهای  $k=\{2,8\}$  می بینیم همگرایی در حالت  $k=2$  کمتر است. زیرا مقدار پارامتر بردار تقویتی بیشتر بوده و پاداش کمتری به عاملها داده می شود. به عبارتی هر چه پارامتر جستجو بیشتر شود سرعت همگرایی بیشتر می شود.

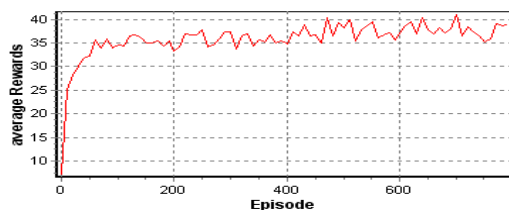
اعمال بازیکنان یعنی  $A2, A1$  بصورت چهار عمل (شمال، جنوب، شرق و غرب یا بالا، پایین، چپ و راست) تعریف می شود. مجموعه فضای حالات بصورت  $S=\{s/s=(l_1, l_2)\}$  تعریف می شود که هر حالت  $s=(l_1, l_2)$  مختصات عاملهای ۱ و ۲ را نشان می دهد. عاملها نمی توانند در یک مختصات یکسان قرار گیرند. اگر دو عامل سعی در حرکت به یک مربع یکسان داشته باشند حرکت هر دو با شکست مواجه می شود. تعداد حالات ممکن در بازی برابر  $8 \times 7 = 56$  است. اگر عاملها به دو مربع مختلف غیر هدف بروند هر دو پاداش صفر را دریافت می کنند و اگر یکی به هدف برسد ۱۰۰ واحد پاداش می گیرد و در صورت برخورد هر دو یک واحد جریمه می شوند و در موقعیت قبلی می مانند. در هر دو بازی گذار از حالتی به

$b$  برابر  $0.002$  برای هر یک از اتوماتاها در نظر گرفته شده است. پارامتر  $K=\{1,3,6\}$  در نظر گرفته شده است. برای نمایش نتایج آزمایشگاه پاداش به دست آمده در هر اپیزود استفاده شده است. ابتدا مقایسه ای بین نتایج تولیدی با نتایج  $Nash-Q$  آورده شده است بطوریکه دیده می شود میانگین پاداش بدست آمده در هر دو الگوریتم تقریباً برابر هستند. با تغییرات پارامتر  $K$  نتایج در لحظه شروع متفاوت و همچنین شیب رسیدن به نقاط تعادل متفاوت میباشد. آزمایشها به خاطر استحکام نتایج ۵۰۰ بار تکرار شده و در هر بار ۷۰۰۰ اپیزود آزمایش شده است. در هر اپیزود جدید هر عامل بطور تصادفی یک مختصات جدید (بجز هدف) را خواهد گرفت. شکل ۶ میانگین پاداش بدست آمده را در هر اپیزود نشان می دهد. همانطور که در شکل دیده میشود در اینجا نیز نقش  $k$  در میزان گرفتن پاداش عاملها مشهود است.

در سری سوم آزمایشها نقش پارامترهای پاداش و جریمه یعنی  $a, b$  مورد بررسی قرار میگیرد. با در نظر گرفتن مقدار  $a=0.2, b=0.002$  می بینیم شیب رسیدن به نقطه تعادل بالا رفته و در فاصله زمانی کمتری به پاداش مطلوبتر می رسد با توجه به تغییر پارامتر  $b=0$  می بینیم اتوماتاهای  $LRI$  رفتار بهتری را از خود نشان می دهند. شکل ۷ نتایج بدست آمده را با پارامترهای مختلف نشان می دهد.



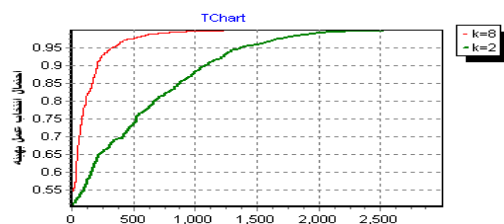
$a=0.2, b=0.005, k=2$



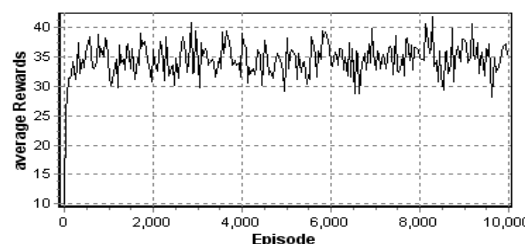
$a=0.2, b=0.0, k=2$

شکل ۷. بررسی رفتار پارامترهای  $a, b$

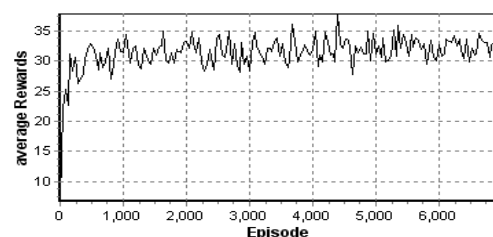
در آزمایشهای نهایی نقش پارامتر  $k$  در بهبود نتایج



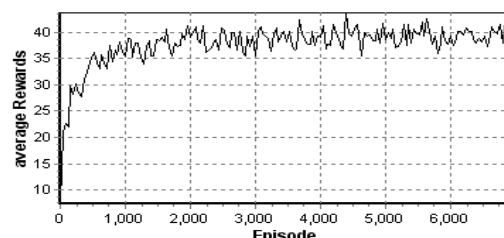
شکل ۵- نمودار تغییرات احتمال انتخاب عمل در خانه شروع به ازاء نرخ یادگیری  $0.02$  در ۳۰۰۰ تکرار برای اتوماتای یادگیرا



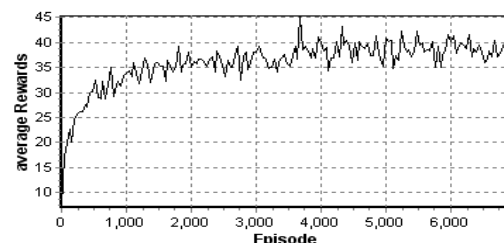
الف ( الگوریتم NASHQ



ب) الگوریتم ارائه شده با  $K=1$



$K=3$



$K=6$

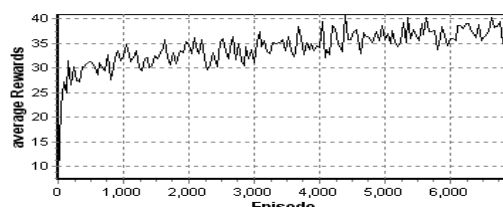
شکل ۶. میانگین پاداش بدست آمده در هر اپیزود در ۵۰۰ آزمایش بازای مقادیر مختلف  $K$  و مقایسه با NASHQ ( $a=0.02, b=0.002$ )

در آزمایشهای سری دوم نیز پارامتر  $a$  برابر  $0.02$  و پارامتر

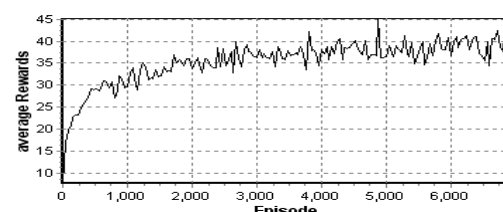
## مراجع

- [1] H. Van Dyke Parunak, "A Practitioners' Review of Industrial Agent Applications, Autonomous Agents and Multi-Agent Systems", v.3 n.4, p.389-407, 2000.
- [2] P. Stone, M. Veloso, "Multiagent systems: A survey from the machine learning perspective," Auton. Robots, vol. 8, no. 3, pp. 345-383, 2000.
- [3] D. Bernstein, S. Zilberstein, N. Immerman, "The complexity of decentralized control of Markov decision processes, Mathematics of Operations Research, Vol. 27, No. 4 pp. 819-840, 2002.
- [4] C. Claus and C. Boutilier, "Sequential optimality and coordination in multiagent systems". In Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence, 1999.
- [5] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," ICML-94, pp. 157-163, 1994.
- [6] J. Hu and M. P. Wellman, "Online Learning about Other Agents in a Dynamic Multiagent System, Journal of Cognitive System Research 2(1), pp. 67-79, Elsevier Science, 2001.
- [7] J. Hu and M. P. Wellman, "Nash Q-Learning for General-Sum Stochastic Games", Journal of Machine Learning Research, 4(Nov):pp. 1039-1069, 2003.
- [8] L. Busnui, R. Babuska, B. Schutter "A Comprehensive Survey of Multiagent Reinforcement Learning". IEEE Transaction on System, Man, Cybern., vol. 38, no.2, pp.156-171, 2008.
- [9] 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.
- [10] M. Song, J. Bai, 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.
- [11] K.S. Narendra and M. A. L. Thathachar, Learning automata: An introduction, Prentice Hall, 1989.
- [12] M.A.L Thathachar and P.S. Sastry, "Varieties of Learning Automata: An Overview", IEEE Transaction on Systems,

بررسی میشود. با توجه به ایده روش بولتزمن در یادگیری  $Q$  ابتدا با پارامتر  $k=1$  نتایج بررسی شده (قابلیت جستجوی بیشتر و به ابتدا عاملها جستجو انجام داده و به مرور زمان با توجه به تغییر پارامتر (افزایش مقدار  $k$ ) امکان بهره برداری از دانسته های عامل میسر شده است. شکل ۸ این نتایج را نشان می دهد.



مقدار اولیه  $k$  برابر ۱ و افزایش 0.0004 در هر مرحله



مقدار اولیه  $k$  برابر ۶ و کاهش 0.0004 در هر مرحله

شکل ۸. بررسی رفتار الگوریتم با تغییر پارامتر  $k$  در طول اجرای الگوریتم.

## ۶-نتیجه گیری

در این مقاله الگوریتمی مبتنی بر اتوماتاهای یادگیر برای حل بازی های مارکوفی ارائه گردید با توجه به الگوریتم ارائه شده و نتایج بدست آمده می بینیم الگوریتم برای حل بازی های مارکوفی با مجموع کلی مناسب است. تعداد تکرار ها، پارامترهای یادگیری و جرائم و نیز مقدار پارامتر  $k$  سرعت رسیدن به تعادل را تعیین می نمایند. تنظیم پارامترهای پاداش و جریمه اتوماتاها می تواند کارایی رسیدن به راه حل بهینه را افزایش داده بطوریکه استفاده از تغییرات پارامتر  $k$  در طول اجرای بازی می تواند زمان رسیدن به راه حل بهینه را بهبود بخشد. با توجه به نتایج به دست آمده اتوماتاهای یادگیر مدل مناسب یادگیری و هماهنگی بین عاملها در سیستمهای چندعامله بوده بطوریکه می تواند به عنوان راه حلی مناسب و کارا در بازی های مارکوفی به کار روند.



Man, and Cybernetics-Part B:  
Cybernetics, Vol. 32, No. pp. 6, 711-722,  
2002.

- [13] B. Jafarpour ,multi agent Cooperatiion  
using LA and PSO ,Msc Thesis,  
Computer Engineering and Information  
Technology Department,Amirkabir  
University of Technology, 2007.

زیر نویس ها

---

- <sup>1</sup> Markov Decision Process (MDP)
- 2 Markov Game
- 3 Partial Observability Markov decision process
- 4Decentralized Partially Observable Markov decision process
- 5 Stochastic Game
- 6 Multi Agent Reinforcement Learning
- 7 Evaluation Function
- 8 Dynamic Programming
- 9 Fully Cooperative
- 10 Nash Equilibrium
- 11 Joint Action and Joint Strategy
- 12 Variable Sturcture Learning Automata
- 13 Multi Learning Automata
- 14 Exploration
- 15 Exploitation