

# همکاری در سیستمهای چند عامله با استفاده از اتوماتاهای یادگیر

محمد رضا خجسته و محمد رضا میبیدی

شناختی خود و نیز براساس تأثیر مستقیم بر ورودی‌های قابل دریافت همدیگر از طریق فعالیتهای ارتباطی با همدیگر هماهنگ می‌کنند. دیگر عاملهای موجود در محیط نیز که اهدافی متقابل با هدف درازمدت این تیم دارند، حریفان یا دشمنان این تیم محسوب می‌شوند.

بعنوان بستری برای تست و پیاده‌سازی سیستمهای چندعامله می‌توان به فوتبال رباتها یا روبوکاپ<sup>۲</sup> اشاره کرد. محیط فوتبال رباتها مثالی از یک محیط پیچیده است که در آن چند عامل باید جهت رسیدن به اهداف تیمی، با هم همکاری کنند [۳] تا [۶]. فرآیندهای رفتاری و تصمیم‌گیری می‌توانند از ساده‌ترین رفتارها، همانند حرکت مستقیم به طرف توپ تا پیچیده‌ترین استدلال‌ها که استراتژی‌های تیم خود و تیم مقابل را روشن می‌سازند، تشکیل شوند.

به دلیل وجود پیچیدگیهای موجود در بسترهایی مانند بستر شبیه‌سازی فوتبال رباتها و رویارویی عاملهای موجود در چنین محیطهایی با حالات بسیار زیاد، متنوع و متغیر، ناگزیر به استفاده از روشهای یادگیری ماشین می‌باشیم. تأکید در این مقاله بر روی سیستمهایی متشکل از چند عامل خودمختار است که می‌توانند در محیطهای زمان واقعی، نویزی<sup>۳</sup>، نیاز به همکاری<sup>۴</sup> و دارای دشمن با اهداف متقابل<sup>۵</sup> عمل کنند [۲].

اتوماتاهای یادگیر بعنوان مدلی برای یادگیری، در محیطی تصادفی عمل نموده و قادر هستند که براساس ورودی‌های دریافت شده از محیط، احتمال انجام عملیات خود را به روز درآورند تا بتوانند از این طریق کارایی خود را بهبود بخشند. یکی از اهداف این مقاله بررسی کارایی اتوماتای یادگیر در همکاری بین عاملهای عضو یک تیم در یک بستر تست شبیه‌سازی فوتبال رباتها می‌باشد. با استفاده از بستر تست شبیه‌سازی فوتبال رباتها به بررسی کارایی اتوماتای یادگیر در همکاری بین عاملهای عضو یک تیم پرداخته شده است. با پیاده‌سازی تیمهای ۲ نفره، ۵ نفره و ۱۱ نفره از عاملها که هر کدام از آنها به یک اتوماتای یادگیر مجهز شده است و مقایسه آنها با یک تیم بدون یادگیری و یا تیمهای یادگیر دیگر، کارایی اتوماتای یادگیر در یادگیری یک کار تیمی جهت دست یافتن به یک هدف مشترک مورد ارزیابی قرار می‌گیرد.

بدلیل وجود تعداد حالات بسیار زیاد در دامنه‌های چندعامله پیچیده، داشتن روشی برای عمومی‌سازی حالات محیطی امری ضروری است انتخاب مناسب چنین روشی، در تعیین حالات و اعمال عامل نقشی تعیین‌کننده دارد. بهمین دلیل در این مقاله به معرفی تکنیک جدیدی بنام "تکنیک بهترین گوشه در مربع حالت" نیز می‌پردازیم. این روش فضای حالات پیوسته و بسیار وسیع عامل را به فضای حالات گسسته و محدود نگاشت می‌کند، ضمن آنکه سعی می‌شود که مواردی مانند فاصله و زاویه که مفاهیمی پیوسته هستند در تعیین حالت عامل منظور شود. ایده نگاشت و عمومی‌سازی فضای دامنه فعالیت عاملهای موجود در یک تیم به یک شبکه از مربعهای حالت برای اولین بار در این مقاله مطرح گردیده است. در کارهای گذشته از جمله [۲]، از مفاهیم پیوسته مانند فاصله، زاویه و

چکیده: عاملها موجودات نرم‌افزاری هستند که بطور پیوسته و خود مختار در یک محیط که برای آن طراحی شده‌اند کار می‌کنند. برای عاملها نیازهایی از قبیل واکنشی بودن نسبت به کنش‌های محیطی، خود مختاری در انتخاب مسیر و ادامه آن، قابلیت تطبیق و یادگیری و ... ضروری به نظر می‌رسد. امروزه مطالعه سیستمهای مبتنی بر عاملها به یک موضوع مهم آکادمیک تبدیل شده است که کاربردهای تجاری و صنعتی فراوانی را نیز در بر دارد. در سیستمهای چندعامله، چندین عامل هوشمند با قابلیت برقراری ارتباط با یکدیگر، جهت رسیدن به مجموعه‌ای از اهداف، با هم همکاری می‌کنند. بدلیل پیچیدگیهای موجود در محیطهای چندعامله پویا و متغیر نیاز به روشهای یادگیری ماشین در چنین محیطهایی احساس می‌شود. اتوماتای یادگیر یک مدل انتزاعی است که تعداد محدودی عمل را می‌تواند انجام دهد. هر عمل انتخاب شده توسط محیطی احتمالی ارزیابی می‌گردد و پاسخی به اتوماتای یادگیر داده می‌شود. اتوماتای یادگیر از این پاسخ استفاده نموده و عمل خود برای مرحله بعد را انتخاب می‌کند. در این مقاله با استفاده از بستر تست شبیه‌سازی فوتبال رباتها به بررسی کارایی اتوماتای یادگیر در همکاری بین عاملهای عضو یک تیم پرداخته شده است. بدلیل وجود تعداد حالات بسیار زیاد در دامنه‌های چندعامله پیچیده، داشتن روشی برای عمومی‌سازی حالات محیطی، امری ضروری است چرا که انتخاب مناسب چنین روشی، در تعیین حالات و اعمال عامل نقشی تعیین‌کننده دارد. در این مقاله همچنین به معرفی و پیاده‌سازی تکنیک "بهترین گوشه در مربع حالت" پرداخته شده است. با استفاده از این روش فضای حالات پیوسته و بسیار وسیع عامل به فضای حالات گسسته و محدود نگاشته می‌شود. کارایی این تکنیک در عمومی‌سازی حالات محیطی در یک دامنه چندعامله همکاری‌گرا مورد بررسی قرار گرفته است.

کلید واژه: اتوماتای یادگیر<sup>۱</sup>، عامل، سیستمهای چندعامله، فوتبال رباتها، همکاری.

## ۱- مقدمه

عامل، موجود خودمختاری است که ویژگیهایی از قبیل اجتماعی بودن، واکنشی بودن و پیش فعال بودن را دارا می‌باشد. عاملها در محیطی زندگی می‌کنند که می‌تواند باز یا بسته باشد و نیز ممکن است که این محیط عاملهایی دیگر را در خود جای داده باشد. هرچند وضعیت‌هایی وجود دارد که یک عامل می‌تواند خودش به تنهایی عمل کند ولی معمول بر آن است که یک عامل با دیگر عاملها ارتباط متقابل داشته باشد. عاملها در جهت دست یافتن به اهداف خود یا جامعه‌ای که در آن زندگی می‌کنند با یکدیگر همکاری می‌کنند. وقتی که یک گروه از عاملها در یک سیستم چندعامله در یک هدف دراز مدت سهیم باشند، آنها تشکیل یک تیم را می‌دهند.

اعضای یک تیم رفتار خود را براساس سازگار کردن فرآیندهای بازنگری شد.

این مقاله در تاریخ ۳ شهریور ماه ۱۳۸۱ دریافت و در تاریخ ۶ مرداد ۱۳۸۲ بازنگری شد.

محمد رضا خجسته، آزمایشگاه محاسبات نرم، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، تهران، ایران (email: khojasteh@ce.aut.ac.ir).

محمد رضا میبیدی، آزمایشگاه محاسبات نرم، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، تهران، ایران (email: meybodi@ce.aut.ac.ir).

1. Learning Automaton

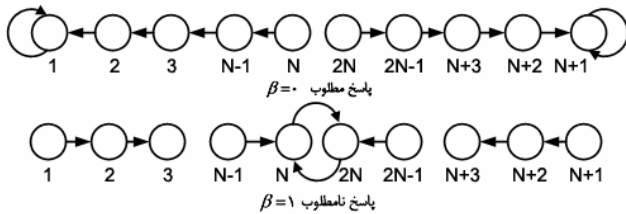
2. Robocup

3. Noisy

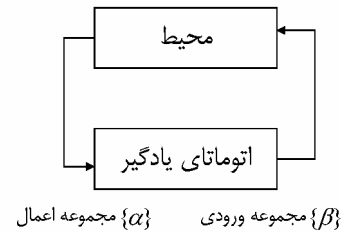
4. Collaborative

5. Adversarial

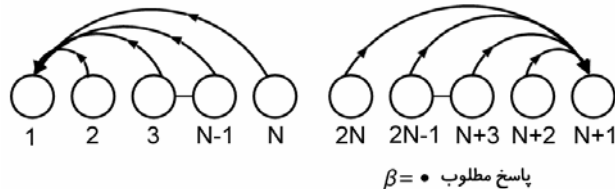
مجموعه ورودی      مجموعه پاسخ



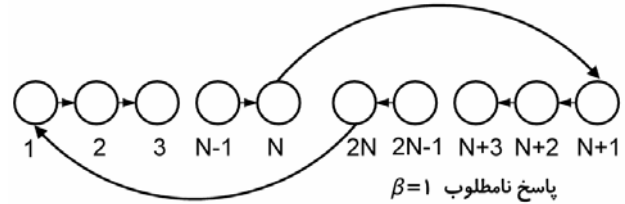
شکل ۲: نمودار تغییر وضعیت اتوماتای  $L_{\gamma, N, \gamma}$ .



شکل ۱: ارتباط بین اتوماتای یادگیر و محیط.



شکل ۴: نمودار تغییر وضعیت اتوماتای Krinsky.



شکل ۳: نمودار تغییر وضعیت اتوماتای  $G_{\gamma, N, \gamma}$ .

ارتباط بین اتوماتای یادگیر و محیط را نشان می‌دهد.

## ۲-۲ اتوماتای یادگیر با ساختار ثابت

اتوماتای یادگیر با ساختار ثابت توسط هتایی  $LA = \{\alpha, \beta, F, G, \phi\}$  نشان داده می‌شود که  $\alpha = \{\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_r\}$  مجموعه اعمال اتوماتا،  $\beta = \{\beta_1, \beta_2, \beta_3, \dots, \beta_r\}$  مجموعه ورودی‌های اتوماتا،  $F: \phi \times \beta \rightarrow \phi$  تابعی که براساس پاسخ محیط، وضعیت جدید را می‌یابد،  $G: \phi \rightarrow \alpha$  تابع خروجی که وضعیت کنونی را به خروجی بعدی می‌نگارد و  $\phi(n) = \{\phi_1, \phi_2, \dots, \phi_k\}$  مجموعه وضعیت‌های داخلی اتوماتا می‌باشند.

در ادامه به چند نمونه از اتوماتاهای یادگیر با ساختار ثابت که در این مقاله از آنها استفاده شده است اشاره شده است.

- اتوماتای  $L_{\gamma, N, \gamma}$ : این اتوماتا تعداد پاداش‌ها و جریمه‌های دریافت شده برای هر عمل را نگهداری کرده و تنها زمانی که تعداد جریمه‌ها بیشتر از پاداش‌ها می‌گردد، عمل دیگر را انتخاب می‌کند. نمودار تغییر وضعیت این اتوماتای مطابق شکل ۲ می‌باشد.

- اتوماتای  $G_{\gamma, N, \gamma}$ : در این اتوماتا بر خلاف اتوماتای  $L_{\gamma, N, \gamma}$ ، عمل  $\alpha_1$  حداقل  $N$  بار انجام گیرد (پس از گرفتن  $N$  جریمه) تا اینکه عمل  $\alpha_1$  دوباره انتخاب شود. گراف تغییر وضعیت این اتوماتا برای پاسخ مطلوب مانند اتوماتای  $L_{\gamma, N, \gamma}$  بوده و برای پاسخ نامطلوب مطابق شکل ۳ می‌باشد.

- اتوماتای Krinsky: این اتوماتا زمانی که پاسخ محیط نامطلوب است، مانند  $L_{\gamma, N, \gamma}$  رفتار می‌کند. اما برای پاسخ مطلوب هر وضعیت  $\phi_i (i=1, 2, \dots, N)$  به وضعیت  $\phi_1$  و هر وضعیت  $\phi (i=N+1, N+2, \dots, 2N)$  به وضعیت  $\phi_{N+1}$  می‌رود. بنابراین همیشه  $N$  پاسخ نامطلوب متوالی لازم است تا اتوماتا عمل خود را عوض کند. نمودار تغییر وضعیت این اتوماتا برای پاسخ نامطلوب مانند اتوماتای  $L_{\gamma, N, \gamma}$  بوده و برای پاسخ مطلوب مطابق شکل ۴ می‌باشد.

- اتوماتای Krylov: در این اتوماتا زمانی که پاسخ محیط مطلوب است، تغییر وضعیت مانند اتوماتای  $L_{\gamma, N, \gamma}$  می‌باشد. اما زمانی که پاسخ محیط نامطلوب باشد، هر وضعیت  $\phi_i (i \neq 1, N, N+1, 2N)$  با احتمال  $0.5$  به وضعیت  $\phi_{i+1}$  و با احتمال  $0.5$  به وضعیت  $\phi_{i-1}$  مطابق شکل ۵ منتقل می‌شود.

نظایر آن برای عمومی‌سازی فضای حالات عامل استفاده شده است. ادامه مقاله بصورت زیر سازماندهی شده است. در بخش ۲ به اختصار به شرح اتوماتاهای یادگیر می‌پردازیم. در بخش ۳ "تکنیک بهترین گوشه در مربع حالت" برای عمومی‌سازی فضای حالات محیطی معرفی می‌گردد. در بخش ۴ به بحث و بررسی موضوع همکاری بین اعضای یک تیم چند عامله می‌پردازیم. در بخش ۵ نتایج آزمایشات ارائه می‌شود. بخش نهایی مقاله نتیجه‌گیری است.

## ۲- اتوماتاهای یادگیر

اتوماتاهای یادگیر مدل‌های انتزاعی هستند که در محیطی تصادفی عمل نموده و قادر هستند که براساس ورودی‌های دریافت شده از محیط، احتمال انجام عملیات خود را به روز درآورده تا بتوانند از این طریق کارایی خود را بهبود بخشند. یک اتوماتای یادگیر تعداد محدودی عمل را می‌تواند انجام دهد. هر عمل انتخاب شده توسط محیطی احتمالی ارزیابی می‌گردد و پاسخی به اتوماتای یادگیر داده می‌شود. اتوماتای یادگیر از این پاسخ استفاده نموده و عمل خود برای مرحله بعد انتخاب می‌کند [۷] و [۸].

اتوماتاهای یادگیر به دو گروه تقسیم می‌گردند:

- الف- اتوماتای یادگیر با ساختار ثابت<sup>۱</sup>
- ب- اتوماتای یادگیر با ساختار متغیر<sup>۲</sup>.

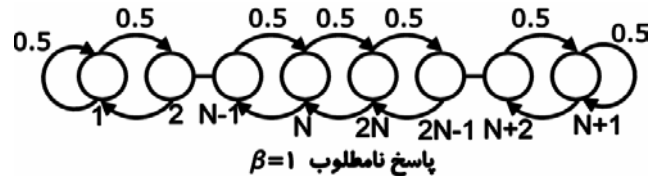
### ۱-۲ محیط

محیط را می‌توان توسط سه تایی  $E \equiv \{\alpha, \beta, c\}$  تعریف نمود که  $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه ورودی‌ها،  $\beta = \{\beta_1, \beta_2, \dots, \beta_r\}$  مجموعه خروجی‌ها و  $c = \{c_1, c_2, \dots, c_r\}$  مجموعه احتمالهای جریمه شدن می‌باشند.

هرگاه  $\beta_i$  دو مقداری باشد، محیط از نوع  $P$  می‌باشد. در چنین محیطی  $\beta_i = 1$  به عنوان جریمه و  $\beta_i = 0$  به عنوان پاداش در نظر گرفته می‌شود.  $c_i$  احتمال این است که عمل  $\alpha_i$  نتیجه نامطلوب داشته باشد می‌باشد. در محیط پایدار<sup>۳</sup> مقادیر  $c_i$  بدون تغییر باقی می‌مانند، حال آن که در محیط ناپایدار<sup>۴</sup> این مقادیر در طی زمان تغییر می‌کنند. شکل ۱

1. Fixed Structure
2. Variable Structure
3. Stationary
4. Non-Stationary

$X-22$	$X-12$	$X-2$	$X+8$	$X+18$
$X-21$	$X-11$	$X-1$	$X+9$	$X+19$
$X-20$	$X-10$	$X$	$X+10$	$X+20$
$X-19$	$X-9$	$X+1$	$X+11$	$X+21$
$X-18$	$X-8$	$X+2$	$X+12$	$X+22$

شکل ۶: ۲۴ مربع اطراف مربع دربردارنده عامل مورد نظر در مربع  $X$ .

شکل ۵: نمودار تغییر وضعیت اتوماتای Krylov.

باشند، هدف پیدا کردن تابع  $f: S \rightarrow V$  می‌باشد. بدین صورت، پس از داشتن تابع  $f$ ، عامل می‌تواند از  $V$  برای یادگیری عمل مناسب خود در محیط استفاده نماید. علاوه بر این تابع، یک تابع تقسیم وظایف  $P: S \rightarrow M$  نیز مورد نیاز است. تابع  $P$  مجموعه حالات دامنه را بین عملهای موجود در محیط تقسیم می‌کند. این تابع فضای حالت را به  $|M|$  بخش مجزا تقسیم می‌کند و هر بخش به حداقل یک عامل (جهت یادگیری و عمل در آن بخش) سپرده می‌شود.

با توجه به توضیحات فوق و با فرض این که مجموعه اعمال عامل  $A$  باشد، عامل در هر کدام از  $|V|$  حالت، دارای  $|A|$  عمل ممکن خواهد بود و بدین ترتیب، مجموعه مورد یادگیری عامل حداکثر شامل  $|A| \times |V|$  عضو خواهد بود. با انتخاب مناسب مجموعه‌های  $V$  و  $A$ ، امکان یادگیری مناسب با مثالهای محدود در دامنه‌ای پیچیده و همزمان فراهم می‌شود. مجموعه‌های  $V$  و  $A$  باید به گونه‌ای انتخاب شوند که تا حد امکان در برگیرنده کلیه حالات و اعمال باشند و نگاشت‌های خوبی از مجموعه‌های حالات و اعمال ممکن در دامنه محیط در بردارنده عامل محسوب شود.

روش پیشنهادی برای عمومی‌سازی محیط بدین صورت است که فضای مستطیل شکل زمین فوتبال که محیط دامنه چندعامله محسوب می‌شود به ۱۵۰ مربع یکسان با ضلع ۷ متر تقسیم می‌شود. از این طریق محیط پیوسته پیرامون عامل به محیطی گسسته تبدیل می‌گردد. در هر لحظه از زمان بازی، هر عامل درون یکی از این مربع‌ها قرار دارد. با توجه به این نکته که هر عامل دارای دیدی محدود می‌باشد برای عاملی که در مربع حالت  $X$  قرار دارد، ۲۴ مربع اطراف آن را بعنوان شعاع دید آن عامل در نظر می‌گیریم. به شکل ۶ توجه کنید.

در این شکل، مربعهای  $X+10$ ،  $X+11$ ،  $X+1$ ،  $X-9$ ،  $X-10$ ،  $X-11$ ،  $X-1$  و  $X+9$  را مربعهای بلافاصله اطراف مربع  $X$  می‌نامیم. شماره‌گذاری مربع‌ها از چپ به راست و بصورت ستونی و با شروع از شماره صفر برای اولین مربع (منتهاالیه چپ و بالا) تا شماره ۱۴۹ برای آخرین مربع (منتهاالیه راست و پایین) انجام گرفته است. هر ستون دارای ۱۰ مربع و هر سطر دارای ۱۵ مربع می‌باشند. در شکل‌های ۷ و ۸، مربع‌های ۸ جهت اطراف مربع دربردارنده عامل (با اختیار جهت شمال به سوی دروازه حریف) نشان داده شده است. در روش پیشنهادی برای هر کدام از ۸ مربع اطراف عامل یعنی برای هر یک از ۸ جهت (شمال غربی، شمال، شمال شرقی، شرق، جنوب شرقی، جنوب، جنوب غربی و غرب) اطراف آن، یک کمیت عددی در نظر گرفته شده است.

به عنوان مثال، اگر عامل واقع در مربع حالت  $X$ ، یکی از عامل‌های حریف را در درون یکی از مربع‌های شمال غرب خود (مربعهای  $X+9$ ،  $X+8$ ،  $X+18$ ،  $X+19$ ) ببیند، یک مقدار عددی منفی به کمیت عددی متناظر با جهت شمال غرب خود اضافه می‌کند. این وضعیت برای عملهای هم تیمی برعکس است. به عنوان مثال، اگر عامل صاحب توپ (واقع در مربع حالت  $X$ ) یکی از عامل‌های خودی را در درون یکی از

### ۳-۲ اتوماتای یادگیر با ساختار متغیر

اتوماتای یادگیر با ساختار متغیر توسط ۴ تائی  $\{\alpha, \beta, p, T\}$  نشان داده می‌شود که در آن  $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه عمل‌های اتوماتا،  $\beta = \{\beta_1, \beta_2, \dots, \beta_m\}$  مجموعه ورودی‌های اتوماتا،  $p = \{p_1, p_2, \dots, p_r\}$  بردار احتمال انتخاب هر یک از عملها و  $p(n+1) = T[\alpha(n), \beta(n), p(n)]$  الگوریتم یادگیری می‌باشد. در این نوع از اتوماتاها، اگر عمل  $\alpha_i$  در مرحله  $n$  انتخاب شود و پاسخ مطلوب از محیط دریافت نماید، احتمال  $p_i(n)$  افزایش یافته و سایر احتمالات کاهش می‌یابند و برای پاسخ نامطلوب احتمال  $p_i(n)$  کاهش یافته و سایر احتمالات افزایش می‌یابند. در هر حال، تغییرات به گونه‌ای صورت می‌گیرد تا حاصل جمع  $p_i(n)$ ‌ها همواره مساوی یک باقی بماند. الگوریتم زیر یک نمونه از الگوریتم‌های یادگیری خطی در اتوماتای با ساختار متغیر است

الف- پاسخ مطلوب

$$p_i(n+1) = p_i(n) + a[1 - p_i(n)]$$

$$\forall j \quad j \neq i \quad p_j(n+1) = (1-a)p_j(n)$$

ب- پاسخ نامطلوب

$$p_i(n+1) = (1-b)p_i(n)$$

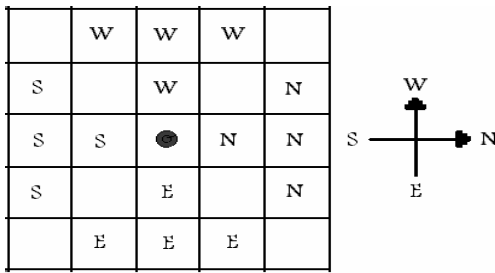
$$p_j(n+1) = \frac{b}{r-1} + (1-b)p_j(n) \quad \forall j \quad j \neq i$$

در روابط فوق،  $a$  پارامتر پاداش و  $b$  پارامتر جریمه می‌باشد. با توجه به مقادیر  $a$  و  $b$  سه حالت زیر را می‌توان در نظر گرفت. زمانی که  $a$  و  $b$  با هم برابر باشند، الگوریتم را  $L_{RP}$  می‌نامیم. وقتی که  $b$  از  $a$  خیلی کوچکتر باشد، الگوریتم را  $L_{REP}$  می‌نامیم و هنگامی که  $b$  مساوی صفر باشد، الگوریتم را  $L_{RI}$  می‌نامیم. حافظه و زمان مورد نیاز برای پیاده‌سازی اتوماتاهای یادگیر با ساختار ثابت  $O(1)$  و برای اتوماتاهای یادگیر با ساختار متغیر  $O(m)$  می‌باشد که  $m$  تعداد اعمال اتوماتا است. برای مطالعه بیشتر درباره اتوماتاهای یادگیر می‌توان به مراجع [۷] تا [۱۱] مراجعه نمود.

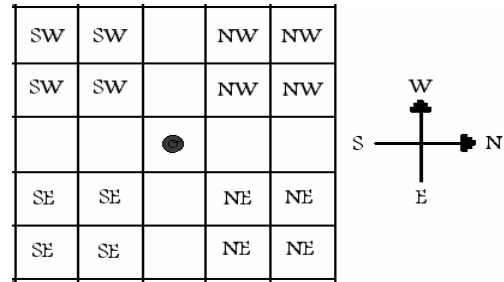
### ۳- تکنیک بهترین گوشه در مربع حالت

همان گونه که قبلاً اشاره شد، تعداد حالات در دامنه فوتبال روبوتیک شبیه‌سازی شده، بسیار زیاد است و لذا امکان در نظر گرفتن کلیه این حالات برای یک عامل عملاً غیرممکن است. بهمین دلیل ایجاد یک روش عمومی‌سازی مناسب از حالات محیطی امری ضروری است. با فرض آنکه مجموعه حالات دامنه  $S$  و مجموعه حالات نگاشت  $V$

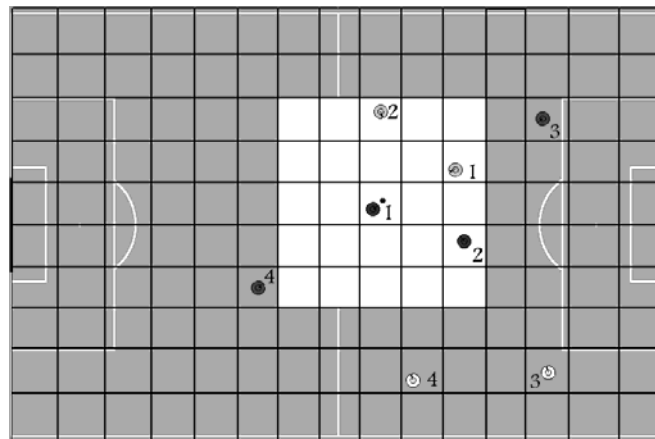
1. Linear Reward Penalty
2. Linear Reward Epsilon Penalty
3. Linear Reward Inaction



شکل ۸: مربعهای شمال، شرق، غرب، و جنوب مربع دربردارنده عامل با توجه به موقعیت عامل مورد نظر و چهار جهت اصلی.



شکل ۷: مربعهای شمال غرب، شمال شرق، جنوب غرب، و جنوب شرق مربع دربردارنده عامل با توجه به موقعیت عامل مورد نظر و چهار جهت.



شکل ۹: وضعیت بازیکنان هم تیمی و حریف با توجه به فضای محلی عامل صاحب توپ در یک بازی ۴ در مقابل ۴.

می‌دهد. یعنی، عامل صاحب توپ در هر سیکل، با توجه به اعداد بدست آمده برای ۸ گوشه مربع حالات در بردارنده خود، "بهترین گوشه در مربع حالت" حالت خود را محاسبه می‌کند. درباره چگونگی انتخاب ۲۴ مربع اطراف هر عامل که فضای محلی یا مهم عامل را تشکیل می‌دهد می‌توان به شکل ۹ مراجعه کرد. در این شکل، فضای محلی و مهم برای یک عامل در روش بهترین گوشه در مربع حالت، در مقایسه با کل فضای بازی نشان داده شده است. این فضا بجای کل فضای بازی برای تصمیم‌گیری محلی یک عامل مورد استفاده قرار می‌گیرد.

شکل ۹ یک بازی ۴ در مقابل ۴ را به تصویر کشیده است. در این شکل، برای بازیکن صاحب توپ، بازیکن ۲ هم تیمی و بازیکنان ۱ و ۲ از تیم حریف صاحب اهمیت هستند (به این دلیل که درون فضای محلی وی قرار دارند) و سایر بازیکنان نقشی در تصمیم‌گیری بازیکن صاحب توپ ندارند.

با توجه به این که حالت هر عامل در محیط در هر لحظه تا حد زیادی به فاصله و زاویه دیگر عامل‌ها نسبت به آن عامل بستگی دارد نداشت بایستی بصورتی انجام گیرد که حالت‌های یکسان به گوشه‌ای واحد و حالت‌های متفاوت به گوشه‌های متفاوت (در مربع حالت عامل) نگاشت شود.

سناریوی کلی برای هر عامل در شبیه‌سازیهای انجام گرفته بدین صورت است که اگر عامل صاحب توپ، بازیکنی غیر از خود را در مسیر به سمت دروازه حریف ببیند، با تعیین حالت خود و انتخاب عمل بهینه در آن حالت (با استفاده از اتوماتای یادگیر وابسته به آن حالت)، سعی می‌کند حرکتی را انجام دهد که در جهت هدف تیمی که همانا برد است باشد.

اگر بازیکنی بازیکن تیم خودی را صاحب توپ ببیند، موقعیت خود و دیگر بازیکنان موجود در حوزه دید خود را به بازیکن صاحب توپ اعلام می‌کند تا بازیکن صاحب توپ، حتی‌المقدور، بهترین تخمین را از حالتی که در آن قرار دارد داشته باشد. این امر با توجه به این نکته که دید هر بازیکن

مربع‌های جنوب غرب خود (مربع‌های  $X-21$  و  $X-22$  و یا  $X-11$  و  $X-12$ ) ببیند، یک مقدار عددی مثبت به کمیت عددی متناظر با جهت جنوب غرب خود اضافه می‌کند. لازم به ذکر است که اندازه این افزایش (منفی و یا مثبت) در ضریب اطمینان هر جهت با معکوس فاصله عامل حریف و یا عامل هم تیمی از عامل مورد نظر (و در واقع صاحب توپ) در آن جهت، متناسب است. در واقع، عامل‌های نزدیک‌تر (و درون فاصله محلی تا ۲۴ مربع حالت پیرامون عامل صاحب توپ) تاثیر بیشتر و عامل‌های دورتر (و درون فاصله محلی تا ۲۴ مربع حالت پیرامون عامل صاحب توپ) تاثیر کمتری را بر روی ضرایب اطمینان عامل صاحب توپ خواهند داشت و عامل‌های خارج از فاصله محلی (خارج از ۲۴ مربع حالت پیرامون عامل)، تاثیری در تغییر ضرایب اطمینان عامل صاحب توپ ندارند.

هر عامل با محاسبه موقعیت تمامی عامل‌های هم تیمی و حریف پیرامون خود (درون ۲۴ مربع)، ۸ عدد در اختیار خواهد داشت. با متناظر کردن این اعداد با ۸ جهت پیرامون عامل صاحب توپ، عامل صاحب توپ دارای ۸ حالت و ۸ عمل (که متناظر با ارسال توپ به سمت مرکز ۸ مربع بلافاصله پیرامون عامل می‌باشند) خواهد بود. لازم بذکر است که تقسیم زمین به مربع‌های  $7 \times 7$  متر مربعی براساس تجربه کسب شده در کار با محیط شبیه‌ساز و میزان جابجایی توپ در اثر هر ضربه حاصل شده است و عملاً وابسته به دامنه انتخابی است. بدین ترتیب، فضای حالات پیرامون عامل به ۸ حالت کاهش می‌یابد و مشکل نگاشت حالات متفاوت به یک حالت بخصوص تا حد زیادی از بین می‌رود.

ممکن است در یک لحظه چندین عامل هم تیمی و یا حریف در درون ۲۴ مربع حالت پیرامون مربع حالت عامل وجود داشته باشند. حال اگر عامل یا عامل‌هایی درون همان مربعی قرار داشته باشند که عامل صاحب توپ در آن قرار دارد، عامل صاحب توپ بر حسب آنکه عامل‌های فوق در کدام یک از ۸ جهت خود قرار دارند کمیت‌های عددی خود را تغییر

جدول ۱: نتایج میانگین گل‌های زده در ۵۰ بازی بین تیمهای یادگیر ۲ نفره با تیم بدون یادگیری ۲ نفره.

	۱	۲	۳	۴	۵	۶
$L_{T,N,T}$	۲-۱	۳-۳	۴-۷	۴-۹	۵-۱۲	۵-۱۶
$G_{T,N,T}$	۱-۲	۱-۵	۲-۸	۲-۱۰	۳-۱۳	۵-۱۶
Krinsky	۱-۲	۲-۴	۳-۷	۳-۱۱	۳-۱۵	۴-۱۷
Krylov	۱-۲	۱-۴	۲-۶	۳-۹	۳-۱۱	۳-۱۲
Q	۱-۳	۲-۵	۲-۱۰	۳-۱۳	۳-۱۵	۴-۱۷

همکاری مورد ارزیابی قرار دهیم.

#### ۴-۱ شبیه‌سازیها برای تیم ۲ نفره

چندین سری شبیه‌سازی انجام گرفته است. در اولین سری شبیه‌سازیها، به پیاده‌سازی تیمهایی ۲ نفره از عاملها می‌پردازیم. این سری از شبیه‌سازیها با دو روش برای تعیین حالت هر عامل در محیط خود انجام میگردد. یکی از این دو روش یک عمومی‌سازی ساده و دیگر روش، تکنیک "بهترین گوشه در مربع حالت" که قبلاً در این مقاله به آن اشاره شده است می‌باشد. در روش عمومی‌سازی ساده، کلیه حالات محیط به ۴ حالت برای بازیکن دارای توپ و ۴ حالت برای بازیکن بدون توپ خلاصه می‌شود [۱۳]. برای هر یک از این چهار حالت یک اتوماتای یادگیر با ساختار ثابت و با عمق حافظه ۳ در نظر گرفته می‌شود. هر اتوماتا، دارای یکی از دو عمل پاس به هم تیمی و یا شوت به طرف دروازه حریف می‌باشد.

نتایج آزمایشها اولیه [۱۳] نشان دادند که تیم دارای اتوماتای یادگیری در مقایسه با یک تیم بدون یادگیری به سرعت یاد می‌گیرد که در چه حالتی، بایستی چه عملی را انجام دهد و به همین دلیل می‌تواند براحتی بر حریف خود غلبه کند. به دلیل این که یادگیری همزمان با بازی انجام می‌گیرد تیم اتوماتای یادگیر می‌تواند در حین بازی خود را با نحوه بازی تیم حریف تا حد زیادی تطبیق دهد. تعداد ۵۰ بازی را بین تیمهای اتوماتای با ساختار ثابت و تیمی مبتنی بر یادگیری Q با تیم بدون یادگیری انجام دادیم که نتایج این بازیها در جدول ۱ آمده است. در این جدول، عدد سمت راست نشان‌دهنده تعداد گل‌های زده توسط تیم یادگیر و عدد سمت چپ نشان‌دهنده تعداد گل‌های زده توسط تیم بدون یادگیری می‌باشد. همانگونه که مشاهده می‌شود کلیه بازیها به سود تیمهای یادگیر به پایان رسیده است. لازم بذکر است که در جدول ۱، ستون ۱ به معنی نتیجه تجمعی بازی از سیکل ۰ تا سیکل ۹۹۹، ستون ۲ به معنی نتیجه تجمعی بازی از سیکل ۱۰۰۰ تا سیکل ۱۹۹۹، ...، و ستون ۶ به معنی نتیجه تجمعی بازی از سیکل ۵۰۰۰ تا آخر بازی (سیکل ۵۹۹۹) می‌باشد. همان گونه که این جدول نشان می‌دهد، تفاوت زیادی بین اتوماتاهای یادگیر مختلف در سری شبیه‌سازیهای انجام شده مشاهده نمی‌شود، هرچند که روشهای یادگیری Q و اتوماتای Krinsky اندکی بهتر از دیگر اتوماتاها عمل کرده است.

همان گونه که جدول ۱ نشان می‌دهد، تیم اتوماتای یادگیر موفق شده است تدریجاً در حین بازی، انجام عمل صحیح پاس و شوت را در ۴ حالت تعریف شده برای آن (در این سری از شبیه‌سازیها) فرا بگیرد. به همین دلیل اکثر گل‌های دریافتی تیم اتوماتای یادگیر در نیمه اول (۳۰۰۰ سیکل اول) و اکثر گل‌های زده شده توسط تیم اتوماتای یادگیر در نیمه دوم بازی (۳۰۰۰ سیکل دوم) به ثمر می‌رسد. باید خاطر نشان کرد که روش یادگیری Q که در در دامنه روبروکاب استفاده شده است دارای محدودیتهایی است. به عنوان مثال، در این دامنه، هر عامل دارای

محدود است و هر بازیکن کنار و پشت خود را نمی‌بیند، برای بازیکن صاحب توپ حیاتی است و لذا در این مورد (و تنها در این مورد) از امکان شنوایی (و آنهم بصورت محدود) استفاده شده است. در واقع در چنین حالتی، بازیکن بدون توپ، مدل دنیای خود را به بازیکن صاحب توپ هم تیمی خود اعلام می‌کند و وی را در تعیین حالت و به تبع آن، تعیین عمل مناسب در آن حالت، یاری می‌کند. بازیکن بدون توپ نیز، به مرکز مربع بلافاصله‌ی (از مربع حالت در بردارنده خود) می‌رود که به منظور دریافت توپ (در صورت لزوم) از بازیکن هم تیمی صاحب توپ، مناسب تشخیص می‌دهد. بدین نحو، در شبیه‌سازیهای انجام گرفته، بازیکن بدون توپ نیز می‌تواند دارای حالات و اعمال خاص (حرکت به سمت یکی از ۸ جهت اطراف خود به منظور دریافت توپ) باشد. یعنی بازیکن می‌تواند "حرکت بدون توپ" انجام دهد. باید توجه داشت که عمل مناسب انتخابی عامل صاحب توپ در هر حالت، الزاماً به معنی پاس و یا شوت نمی‌باشد و در بعضی موارد، ممکن است دریبل و یا نگهداشتن موقت توپ به منظور یافتن روزنه‌ای برای انجام عمل مناسب باشد.

همچنین در روش پیشنهادی، بازیکن خود را ملزم به حرکت رو به جلو (در همه حالات) نمی‌داند و در صورت لزوم از مفاهیم اوت کردن توپ، پاس به عقب، و حتی بازکردن بازی و ارسال توپ به فضای خالی (با احتمال دریافت مناسب توپ برای بازیکن هم تیمی و یا خود در چند سیکل بعد) نیز استفاده می‌کند [۱۳].

لازم به ذکر است که یادگیری در کلیه شبیه‌سازیهای انجام گرفته کاملاً چندعامله و توزیع شده است و برخلاف روشهای موجود، بازیکن صاحب توپ (بدون توپ)، قصد خود را از ارسال (دریافت) توپ به (از) بازیکن بدون توپ (صاحب توپ) اعلام نمی‌دارد و بدین ترتیب، هر عامل در انتخاب عمل بهینه خود کاملاً خودمختار است و درحین انجام بازی بدون ایجاد ارتباط با دیگر بازیکنان سعی می‌کنند که با آنها همکاری کنند. شبیه‌سازیهای انجام شده برای تیمهای ۲ نفره، ۵ نفره و ۱۱ نفره [۱۳] که بعداً در این مقاله ارایه خواهد گردید نشان دهنده این قابلیت می‌باشد.

#### ۴- همکاری بین اعضای یک تیم چند عامله با استفاده از اتوماتای یادگیر

تا بحال از روشهای مختلفی از جمله یادگیری Q، الگوریتمهای ژنتیک، درختهای تصمیم‌گیری و یادگیریهای رفتاری [۲] برای یادگیری عامل‌های فوتبالیست استفاده شده است. در این بخش از مقاله، توانایی اتوماتای یادگیر بمنظور ایجاد همکاری بین عامل‌های موجود در یک تیم روبروکاب برای رسیدن به یک هدف مشخص تیمی که همانا برد باشد را مورد بررسی قرار می‌دهیم. در آزمایشهای انجام شده در این مقاله، چند بازیکن فوتبال که هرکدام مجهز به یک اتوماتای یادگیر می‌باشند را در مقابل چند بازیکن بدون توانایی یادگیری و یا دارای روشهای دیگر یادگیری قرار داده‌ایم تا از این طریق مدل اتوماتای یادگیر را در ایجاد

جدول ۲: نسبت گل زده به گل خورده در روشهای مختلف یادگیری با استفاده از تکنیک بهترین گوشه در مربع حالت و در بازیهای ۲ در مقابل ۲ (در حین آموزش)

نسبت گل زده به گل خورده در حین آموزش	نام روش
۱/۲۶	اتوماتای $L_{\gamma N, \gamma}$ با عمق حافظه ۳ (۱۰ بازی)
۱/۱۷	اتوماتای $L_{\gamma N, \gamma}$ با عمق حافظه ۲ (۱۰ بازی)
۱/۲۵	اتوماتای $G_{\gamma N, \gamma}$ با عمق حافظه ۲ (۱۰ بازی)
۱/۳۰	اتوماتای Krylov با عمق حافظه ۲ (۱۰ بازی)
۰/۸۲	اتوماتای Krinsky با عمق حافظه ۲ (۱۰ بازی)
۱/۱۸	اتوماتای $L_{RP}$ (۳۰ بازی)
۱/۰۵	اتوماتای $L_{RP}$ (۱۵ بازی)
۱/۰۷	اتوماتای $L_{REP}$ با نسبت پارامتر ۰/۱ (۱۵ بازی)
۱/۱۲	اتوماتای $L_{REP}$ با نسبت پارامتر ۰/۰۱ (۱۵ بازی)
۱/۰۲	یادگیری Q (۱۰ بازی)

می‌باشد. لهذا برای عمومی‌سازی بهتر حالات محیطی از تکنیک "بهترین گوشه در مربع حالت" استفاده خواهیم کرد. همچنین تعداد اعمال هر عامل که به پاس دادن به هم تیمی و یا شوت به دروازه حریف خلاصه می‌شد نیز افزایش داده شد. کلیه شبیه‌سازیهایی که در ادامه این مقاله آمده است با توجه به دو مورد فوق انجام گرفته است و شبیه‌سازیهای ۲ بازیکن در مقابل ۲ بازیکن مجدداً انجام گرفت.

برای آزمایشهای جدید دو تیم با نامهای "یادگیری کامل" و "نگاشت ثابت" ساخته شدند. در تیم "یادگیری کامل" در هر مربع حالت از فضای محیط عامل، ۸ اتوماتا تعبیه کرده‌ایم. در این تیم یک عامل، اتوماتا و یا اتوماتاهایی که به آنها تجهیز شده است با حرکت در محیط با خود جابجا نمی‌کند. تیم "نگاشت ثابت" یک تیم بهینه شده با دست است. این دو با تیمهای دیگر مقایسه گردیده است [۱۳].

در یک سری از آزمایشها، تیم یادگیر  $L_{RP}$  با استفاده از تکنیک "بهترین گوشه در مربع حالت" را در چند بازی آموزشی در مقابل تیم استفاده می‌کرد، قرار دادیم. مجموعاً ۳۰ بازی بین تیم  $L_{RP}$  با تیم  $L_{\gamma N, \gamma}$  برگزار کردیم که نتیجه ۶ بازی اول آنها در زیر آمده است. در مجموع ۳۰ بازی، ۲۸ بار تیم  $L_{RP}$  برنده بازی بود و یک بار نیز بازی به تساوی کشیده شد [۱۳].

می‌توان از تعداد گل‌های زده و خورده تیمهای یادگیر مختلف زمانی که در مقابل تیم بدون یادگیری بازی می‌کنند استفاده کرد و تاثیر روشهای یادگیری متفاوت را در کارایی تیم زمانی که از تکنیک بهترین گوشه در مربع حالت استفاده می‌شود را مورد بررسی قرار داد. نتایج حاصله که در جدول ۲ آمده است نشان می‌دهد که عامل در تعیین حالت و انتخاب عمل مناسب در آن حالت بهتر عمل می‌نماید. همچنین می‌توان این نتیجه را گرفت که روش بهترین گوشه در مربع حالت، عمومی‌سازی مناسبتری از حالات محیط و نگاشت حالات یکسان به یک حالت واحد را انجام می‌دهد.

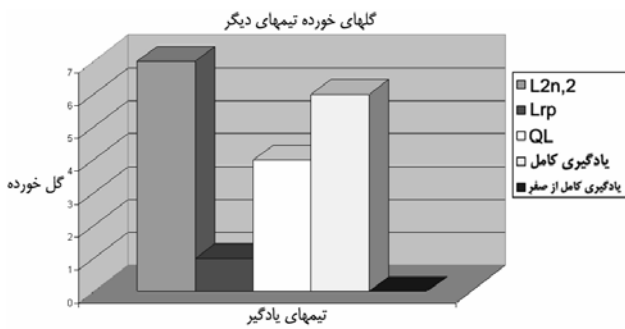
#### ۴-۲ شبیه‌سازیها برای تیم ۵ نفره

قبل از این که به ارائه نتایج شبیه‌سازیها برای تیم ۵ نفره بپردازیم بینیم در چه مواقعی اتوماتای یادگیری که عامل به آن مجهز است پاداش یا جریمه دریافت می‌کند. چگونگی دادن پاداش و یا جریمه به یک عمل به صورت زیر انجام می‌گیرد: اگر یک بازیکن عملی را انجام دهد در

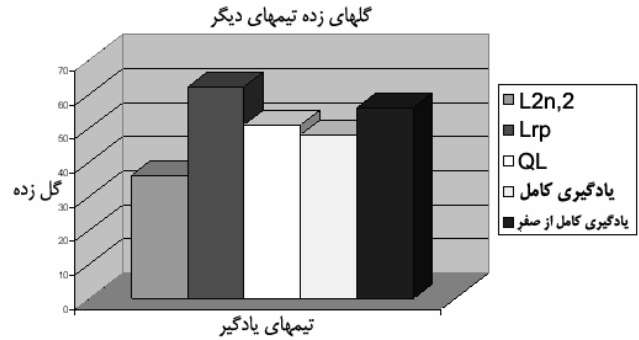
دید محدود است و لذا نمی‌تواند از بقیه عاملها، محیط و همین طور تاثیر (بخصوص درازمدت) اعمال خود بر محیط و دیگر عاملها دیدی کامل داشته باشد. از این رو برای پیاده‌سازی یادگیری  $Q$ ، شبیه‌سازی ساده‌ای از روش یادگیری  $Q$  براساس روش بکار برده شده در [۲] (روش TPOT\_RL) انجام شده است و تنها تاثیرات کوتاه مدت اعمال انجام گرفته توسط عامل، در تغییر مقادیر  $Q$  نقش دارند.

درباره شبیه‌سازیها دو نکته قابل ذکر است. نکته اول این که در کلیه شبیه‌سازیهای انجام شده، از عامل ثالثی به منظور مشاهده روند کار عاملها استفاده نشده است بلکه از خود عامل جهت قضاوت تاثیر عمل خود استفاده کرده‌ایم که این خود می‌تواند دلیلی بر تاکید بر یادگیری در یک سیستم چندعامله در شبیه‌سازیهای انجام شده در این مقاله باشد. نکته دوم این است که عامل نمی‌تواند همه تغییرات در محیط خود را مشاهده کند و همواره حالات مخفی و پوشیده در محیط وجود دارد و به همین دلیل در اکثر موارد، اولین تاثیر قابل مشاهده عمل برای عامل انجام‌دهنده آن عمل، تخمین مناسبی از تاثیر درازمدت عمل انجام شده توسط عامل می‌باشد. قضاوت براساس اولین تاثیر قابل مشاهده توسط (دید محدود) خود بازیکن، در کلیه شبیه‌سازیهای انجام شده در این مقاله رعایت شده است.

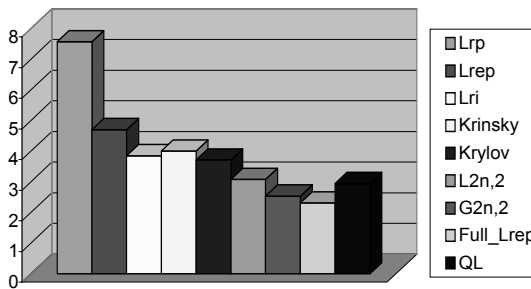
شبیه‌سازیهای دیگر نشان دادند که تیمهای مجهز به اتوماتای یادگیر ۲ نفره بر تیمهای بدون یادگیری با تعداد نفرات بیشتر (۳ یا ۴ نفره) نیز غلبه نمودند. نکته قابل توجه دیگر که از آزمایشها حاصل شد این است که یادگیری چندعامله تا حدی وابسته به حریف است. یعنی مقدار از قبل تعیین شده‌ای برای همگرایی مقادیر حافظه بازیکنان وجود ندارد. این بدان معنی است که بازیکنان در بعضی موارد که از نظر ما منطقاً باید پاس بدهند، یاد گرفته‌اند که بهتر است به سمت دروازه حریف شوت کنند و این یکی از محاسن یادگیری (بخصوص یادگیری همزمان با بازی) می‌باشد. آزمایشهایی که تا بحال انجام گرفته است دارای دو ایراد می‌باشد. اولین ایراد این است که امکان نگاشت حالات بسیار متفاوتی از دنیای پیرامون عامل صاحب توپ به یک حالت واحد وجود دارد و لذا ممکن است عملی که عامل بعنوان عمل صحیح در آن حالت انتخاب کرده است، بهترین عمل ممکن برای آن عامل نباشد. شبیه‌سازیها نشان دادند که عامل در شناخت حالت خود دچار مشکل می‌باشد. برای حل این مشکل می‌توان از روش "عمومی‌سازی" بهتری از حالات محیطی برای عامل استفاده کرد. ایراد دیگر تعداد کم اعمال تعریف شده برای هر عامل



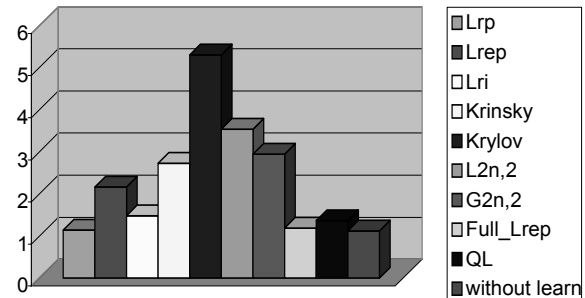
شکل ۱۱: گل‌های خورده تیمهای یادگیر ۵ نفره در ۱۰ شبیه‌سازی (در حین آموزش).



شکل ۱۰: گل‌های زده تیمهای یادگیر ۵ نفره در ۱۰ شبیه‌سازی (در حین آموزش).



شکل ۱۳: مقایسه نسبت میانگین گل‌های زده به خورده در هر بازی پس از آموزش (آزمایشی) در بازیهای تیم بدون یادگیری با تیمهای دیگر.



شکل ۱۲: مقایسه نسبت میانگین گل‌های زده به خورده در هر بازی آموزشی در بازیهای تیم بدون یادگیری با تیمهای دیگر.

بازیهای ۲ نفره در هیچ دیداری بازنده نبودند. نتایج آزمایشها در اشکال ۱۰ و ۱۱ آمده است.

### ۴-۳ شبیه‌سازیها برای تیم ۱۱ نفره

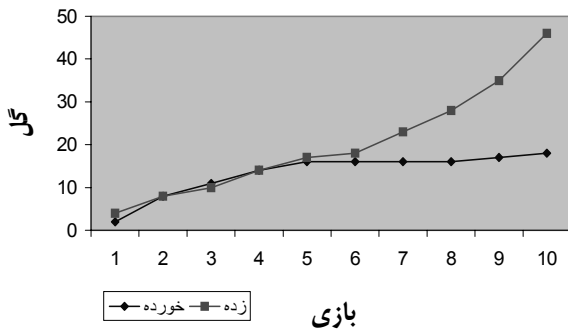
شبیه‌سازیهای ارایه شده در این قسمت به بررسی یادگیری در نحوه همکاری در بین عاملهای موجود در یک تیم ۱۱ نفره و مقایسه آن با دیگر تیمها می‌پردازد. با توجه به این که در بازیهای ۱۱ نفره در مقابل ۱۱ نفر، ممکن است تعداد گل‌های رد و بدل شده خیلی زیاد نباشند به منظور بررسی نتایج شبیه‌سازیها در این سری (و سریهای بعد)، از معیارهای دیگری هم بمنظور نشان دادن کارایی استفاده می‌شود. پس از بررسی معیارهای مطرح در همکاری بین بازیکنان یک تیم فوتبال، معیارهای زیر بنظر مناسب می‌باشد.

- درصد مالکیت توپ توسط تیم خودی در مقایسه با مورد مشابه در تیم حریف در حین بازی.
- درصد گردش توپ در ۱/۳ زمین خودی، ۱/۳ میانی زمین، و ۱/۳ زمین حریف در حین بازی.
- حداکثر زمان در اختیار داشتن توپ بصورت ممتد توسط تیم که بر حسب سیکل سنجیده می‌شود.
- حداکثر تعداد رد و بدل‌های متوالی توپ بدون برخورد با حریف.
- میانگین درصد خطای (اعمال) بازیکنان تیم خودی در حین بازی.

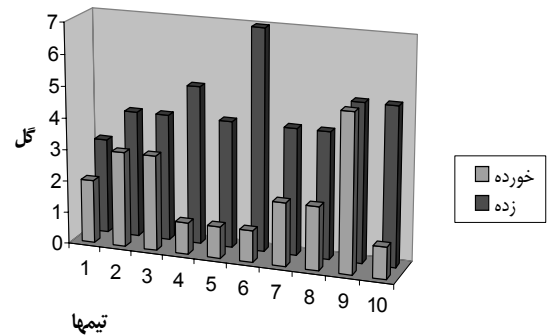
برای سازماندهی ۱۱ بازیکن درون زمین برای هر تیم آرایش تیمی ۳-۳-۴ انتخاب گردیده است. همانند شبیه‌سازیهای قبلی، در این سری از آزمایشها نیز برای مقایسه روند یادگیری در تیمهای یادگیر، یک تیم بدون یادگیری ایجاد کردیم. این تیم بدون یادگیری، به غیر از مساله یادگیری در همه موارد از جمله شکل تیمی مشابه تیمهای یادگیر است. در شبیه‌سازیهای انجام شده تیمهای اتوماتای یادگیر موفق شدند پس از انجام تعدادی محدود بازی آموزشی بر حریف بدون یادگیری غلبه کنند. آنها همچنین توانستند تیم "نگاشت ثابت" را شکست دهند. شکل‌های ۱۲

صورتی پاداش می‌گیرد که تیم صاحب توپ بعد از خاتمه این عمل هنوز تیم اوست و توپ به سود تیم او جلوتر رفته است (بطور مثال تیم او گل زده است) و در غیر این صورت بازیکن جریمه می‌شود. همچنین اگر انجام دهنده آخرین عمل، خود من بوده‌ام (یا اگر هنوز نتیجه آخرین عمل خود را ندیده‌ام) و تیم صاحب توپ، تیم حریف است و اگر فاصله موقعیت قبلی توپ تا موقعیت فعلی توپ خیلی زیاد نیست (بطور مثال به تیم حریف گل زده‌ام) بازیکن جریمه می‌شود و در غیر این صورت به او پاداش داده می‌شود. نتیجه عمل که پاداش یا جریمه می‌باشد توسط اتوماتای یادگیر استفاده می‌شود تا حالت داخلی خود را به روز نماید. چگونگی به روز در آمدن حالت اتوماتا بستگی به نوع اتوماتای یادگیر دارد. برای تیمهای ۱۱ نفره که در ادامه این قسمت به آن پرداخته می‌شود روش دادن پاداش یا جزا نیز به همین ترتیب خواهد بود.

در ادامه به بررسی شبیه‌سازیها برای تیمهای ۵ نفره می‌پردازیم. مهمترین تغییری که در این شبیه‌سازیها نسبت به شبیه‌سازیهای قسمتهای قبل انجام گرفته است دادن شکل تیمی به تیم و استفاده از بازیکنان تیم در پستهای تخصصی بوده است. دادن شکل تیمی با توجه به افزایش نفرت تیم جهت برقراری نظم و همکاری هر چه بهتر در بین عاملها بنظر ضروری می‌رسید. علت انتخاب ۵ نفر برای یک تیم، شبیه‌سازی مواردی چون مسابقات رابوکاپ با روباتهایی با اندازه‌های متوسط و بزرگ و نیز فوتبال داخل سالنی بوده است. در شبیه‌سازیهای ۵ نفره، تیمهای یادگیر  $L_{2n,2}$  بعنوان نماینده تیمهای یادگیر با ساختار ثابت،  $L_{RP}$  بعنوان نماینده تیمهای یادگیر با ساختار متغیر،  $Q$  بعنوان نماینده‌ای از تیمهای دارای روش یادگیری به غیر از اتوماتای یادگیر، و "تیم یادگیری کامل" [۱۳] استفاده کردیم. در شبیه‌سازیهای این قسمت، از ۱ دروازه‌بان، ۲ بازیکن کناری چپ و راست، یک دفاع عقب و یک بازیکن حمله استفاده شد است. به هر کدام از این بازیکنان فضای مشخصی از زمین که در واقع محدوده‌ای از مربعهای حالت می‌باشد تخصیص داده شده است. در بازیهای انجام شده، تیمهای یادگیر، همانند



شکل ۱۵: مقایسه گل‌های زده و خورده در طی ۱۰ بازی آموزشی متوالی با پارامترهای تصادفی اضافه شده زمانیکه تیم اتوماتای  $L_{REP}$  با تیم بدون یادگیری بازی می‌کند.



شکل ۱۴: مقایسه گل‌های خورده (میله های جلو) و گل‌های زده (میله های عقب) تیم‌های مختلف در مقابل تیم دستنویس (ثابت)، از چپ به راست: بدون یادگیری، Q، کامل  $L_{REP}$ ،  $L_{RP}$ ،  $L_{RI}$ ،  $L_{REP}$ ،  $L_{RN}$ ،  $Krylov$ ،  $Krinsky$ ،  $G_{RN}$  و  $G_{RN}$ .

تیم Saloo 2001 توسط نودا نوشته شده است. نودا در شرح تیم خود تقلید را بعنوان اولین قدم برای سازگاری یک عامل با عامل‌های دیگر در محیطی چندعامله می‌داند [۱۶]. او از یک شبکه عصبی بازگشتی برای یادگیری و برنامه‌ریزی استفاده نموده است. این شبکه عصبی دارای ۲ وظیفه است: پیشگویی محیط و شناخت بازی. شبکه عصبی استفاده شده می‌تواند بعنوان یک تمیزدهنده بین انواع بازیها آموزش داده شده مورد استفاده قرار گیرد. وی نتیجه‌گیری کرده است که این شبکه به همراه معماری وی می‌تواند انواع گوناگونی از بازیها را مجدداً تولید کند [۱۶].

در آزمایشهای این قسمت از تیم اتوماتای  $L_{REP}$  و تیم یادگیر کامل  $L_{REP}$  استفاده شده است. این انتخاب بدین دلیل بوده است که این تیمها در شبیه‌سازیهای قبلی [۱۲] تا [۱۵] نتایج خوبی را تولید کرده است. در اولین سری از شبیه‌سازیها به ارزیابی تیم اتوماتای یادگیر در شرایطی که نویز در محیط وجود داشته باشد می‌پردازیم. پارامترهای متعددی در کارگزار شبیه‌سازی فوتبال وجود دارند که می‌توان با تغییر آنها، شرایط حاکم بر زمین بازی، حرکت بازیکنان، حرکت توپ و ... را تغییر داد و یا در آنها اختلال ایجاد کرد.

در اولین نمونه از شبیه‌سازیها، اثر پارامترهای rand که معرف مقدار نویز در موارد مختلف هستند مورد بررسی قرار داده شده است. بدینصورت که پارامتر player\_rand را از مقدار ۰/۱ تا ۰/۲، پارامتر ball\_rand را از ۰/۵ تا ۰/۱، و بالاخره پارامتر kick\_rand را از ۰ تا ۰/۱ تغییر دادیم. پارامتر اول باعث ایجاد اختلال (نویز) توسط کارگزار در حرکت بازیکن می‌شود و پارامترهای دوم و سوم، همین نقش را در مورد حرکت توپ در زمین و زدن ضربه به توپ بر عهده دارند. ۱۰ بازی متوالی بین تیم  $L_{REP}$  (با یادگیری از صفر) و تیم بدون یادگیری برگزار نمودیم. هدف از این شبیه‌سازی این بود که ببینیم آیا روشهای یادگیری پیشنهادی در این مقاله می‌تواند در این گونه شرایط نیز موثر واقع شوند. شکل ۱۵ نتایج را (بصورت تجمعی) نشان می‌دهد. در این شکل، کارایی روش یادگیری پیشنهادی (اتوماتای یادگیر  $L_{REP}$ ) به صورت روشن مشاهده می‌شود. همانگونه که شکل نشان می‌دهد، با افزایش تعداد بازیها، تیم اتوماتای یادگیر موفق می‌شود که رفته رفته با شرایط محیطی سازگار شود و فاصله خود را با تیم بدون یادگیری بیشتر و بیشتر نماید. حتی در شرایط با نویز بالا نیز تیم اتوماتای یادگیر می‌تواند خود را با شرایط محیطی سازگار کند. گر چه بدلیل تاکید بر خصوصیت خودمختاری عامل‌های تیمهای پیشنهادی از حداقل ارتباطات بین بازیکنان در حین بازی سود می‌برد ولی بدلیل اینکه کانال ارتباطی بین بازیکنان شلوغ و فاقد اطمینان است، برای همین حجم اندک ارسال (و نه تبادل) اطلاعات اتکا به ارتباطات بین عامل‌ها در چنین دامنه‌هایی قابل اطمینان نیست. تیمهای اتوماتاهای

و ۱۳ نتایج کلی شبیه‌سازیها برای بازیهای انجام شده بین تیمهای یادگیر و تیم بدون یادگیری را در دوسری بازیهای آموزشی (۱۵ بازی) و بازیهای آزمایشی (۳ بازی پس از بازیهای آموزشی) خلاصه می‌کنند. همانگونه که مشاهده می‌شود با تعدادی محدود بازیهای آموزشی تیمهای اتوماتای یادگیر با ساختار ثابت به برتری نسبی می‌رسند. در مورد بازیهای آزمایشی (پس از آموزش و تنها با استخراج مقادیر یادگرفته شده) به نظر می‌رسد که برتری نسبی (در غلبه بر تیم بدون یادگیری) با تیمهای اتوماتای یادگیر با ساختار متغیر می‌باشد. جزئیات این آزمایشها در [۱۳] آمده است. شکل ۱۴ گل‌های زده و خورده تیمهای یادگیر را در مقابل تیم "نگاشت ثابت" خلاصه می‌کند.

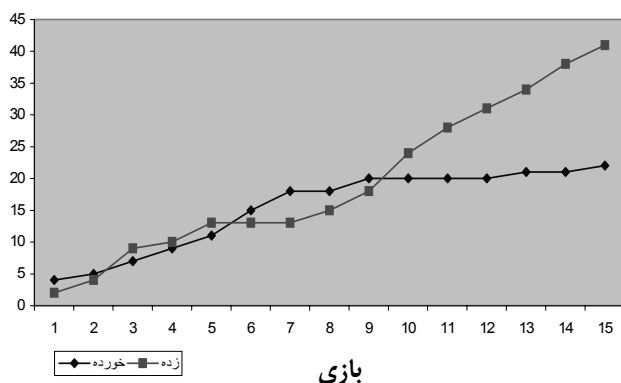
## ۵- تستهای ارزیابی برای تیمهای مبتنی بر اتوماتای یادگیر

در بخشهای قبلی، عملکرد اتوماتاهای یادگیر را در انجام یک کار گروهی در بین عاملهای عضو یک تیم مورد بررسی قرار دادیم. تمام شبیه‌سازیهای انجام شده در موارد فوق، در شرایط طبیعی دامنه برگزار شدند. در این بخش به ارزیابی تیمهای یادگیر از طریق بازی با بعضی از تیمهای شرکت‌کننده در مسابقات رابوکاپ جهانی می‌پردازیم.

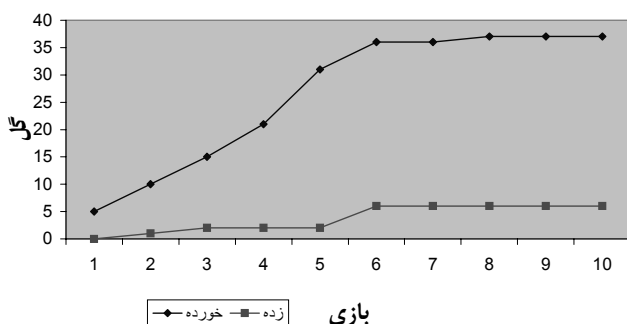
گر چه تیم اتوماتای یادگیر بدلیل پیاده‌سازی نکردن مواردی مانند استراتژی حمله، استراتژی دفاع، مدل کردن حریف، تغییر استراتژی بازی در زمانهای خاص، مربی و ... یک تیم مسابقه محسوب نمی‌شود ولی بازی با تیمهای شناخته شده که در مسابقات رابوکاپ شرکت می‌کنند، می‌تواند ما را در ارزیابی مدل اتوماتاهای یادگیر در یادگیری عامل‌ها یاری کند. لازم به ذکر است که تیم اتوماتای یادگیر یک تیم تحقیقاتی می‌باشد و صرفاً برای مطالعه و بررسی تواناییهای اتوماتاهای یادگیر در همکاری بین عامل‌ها در انجام یک کار گروهی استفاده شده است.

یادآوری می‌شود که کد پایه که برای پیاده‌سازی تیمهای اتوماتاهای یادگیر مورد استفاده قرار گرفته است، کد CMUnited98 [۲] می‌باشد. لازم به ذکر است که این تیم در مقایسه با تیمهای جدیدتر از مهارتهایی فردی و گروهی بالایی بر خوردار نیست. تیمهای Sharif Arvand 2000، Yberoos 2000، FuzzyFoo 2001 و Saloo 2001 که از نظر فردی اختلاف زیادی با تیم ما نداشتند (هر چند بهتر از تیم ما بودند) می‌توانند بعنوان حریف بازیها انتخاب شوند. از بین چهار تیم فوق، تیم Saloo 2001 بدلیل نزدیکتر بودن مهارتهای فردی با تیم خودی مناسب‌ترین تیم برای مسابقه به منظور ارزیابی روش اتوماتای یادگیر تشخیص داده شد و به همین دلیل بازیها با این تیم انجام گرفت [۱۳]. برای مشاهده نتایج بازیها با دیگر تیمهای فوق می‌توان به [۱۳] مراجعه کرد.

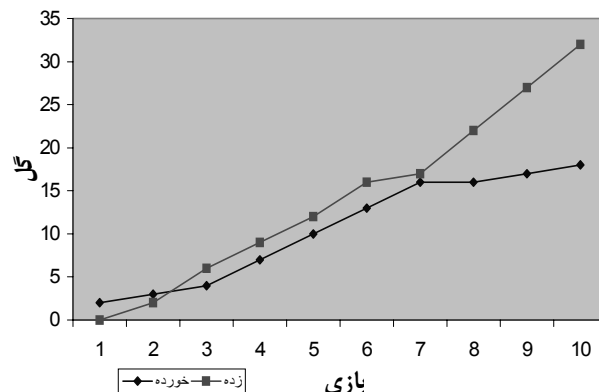




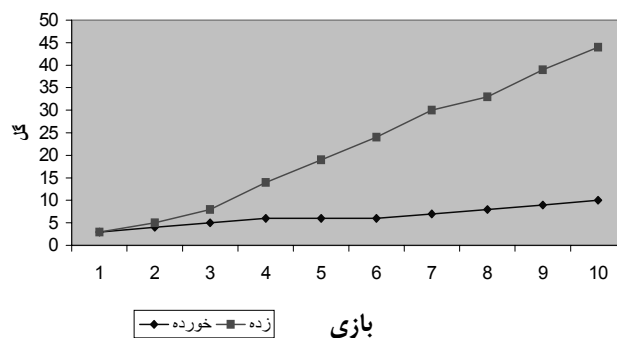
شکل ۱۷: مقایسه گل‌های زده و خورده در طی ۱۰ بازی آموزشی متوالی با برداشتن ۳ بازیکن جناح چپ از تیم یادگیر زمانیکه تیم اتوماتای  $L_{REP}$  با تیم بدون یادگیری بازی می‌کند.



شکل ۱۹: مقایسه گل‌های زده و خورده در طی ۱۰ بازی آموزشی متوالی بین تیم اتوماتای  $L_{REP}$  با تیم بدون یادگیری با افزودن بردار باد (۰ و ۵۰) به زمین بازی (جهت باد، عمود بر عرض زمین بازی و تیم یادگیر در خلاف جهت باد).



شکل ۱۶: مقایسه گل‌های زده و خورده در طی ۱۰ بازی آموزشی متوالی بدون امکان ارتباط مابین بازیکنان بین تیم اتوماتای  $L_{REP}$  زمانیکه با تیم بدون یادگیری بازی می‌کند.



شکل ۱۸: مقایسه گل‌های زده و خورده در طی ۱۰ بازی آموزشی متوالی بین تیم اتوماتای  $L_{REP}$  با تیم بدون یادگیری با افزودن بردار باد (۰ و ۵۰) به زمین بازی (جهت باد، عمود بر طول زمین بازی).

در شبیه‌سازیهای سری بعد، تیم اتوماتای  $L_{REP}$  را (با یادگیری از صفر) در چند بازی در مقابل تیم Saloo قرار دادیم. همانگونه که نتایج نشان می‌دهند [۱۳]، تیم Saloo در تمام بازیهای اولیه برنده بازی بوده است و با تفاضل گل نسبتاً بالایی (میانگین گل زده ۵/۷ گل در مقابل گل خورده ۰/۳ در هر بازی) تیم اتوماتای یادگیر را شکست داده است. آمار بدست آمده از ۷ بازی اولیه در جدول ۳ آمده است. همانگونه که آمار بازی نشان می‌دهد، تیم Saloo برتری مطلق بر تیم اتوماتای یادگیر (در شروع یادگیری خود) دارد.

در ادامه آزمایشها تعداد بازیهای بین تیمها را افزایش دادیم. در طی ۱۵۰ بازی آموزشی (۲۵ ساعت بازی آموزشی متوالی) بین تیم اتوماتای یادگیر با تیم Saloo (که نتایج ۷ بازی اول آنها در جدول ۳ آمده است)، تیم اتوماتای یادگیر رفته رفته بازی بهتری ارائه داد و نتایج بازیها به مساوی و در نهایت به برد پیوسته تیم اتوماتای یادگیر (در بازیهای آخر) انجامید. آمار ۷ بازی پس از آموزش ۲۵ ساعته انجام شده بین تیم اتوماتای یادگیر و تیم Saloo در جدول ۴ آمده است. در این ۷ بازی، تیم اتوماتای یادگیر میانگین گل زده ۳/۶ گل در مقابل گل خورده ۰/۱ در هر بازی را داشته است. همانگونه که از روی آمار بازی نیز مشخص است، تیم Saloo هر چند در تعدادی از معیارها، دارای نتایج بهتری از تیم اتوماتای یادگیر می‌باشد، اما برتری مطلق خود را بر تیم اتوماتای یادگیر آن گونه که در شروع بازی وجود داشت را از دست داده است.

یکی دیگر از معیارهایی که برای ارزیابی کارایی همکاری در بین عاملهای تعریف شده است، معیار درصد اعمال صحیح انجام گرفته و همچنین خطای یک بازیکن در طول بازی است. در این قسمت و به منظور مقایسه، میانگین درصد خطای بازیکنان تیم یادگیر را در ۷ بازی

یادگیر حتی در چنین دامنه‌هایی می‌توانند موفق باشند. بدین منظور، ۱۰ بازی متوالی آموزشی بین تیم اتوماتای یادگیر  $L_{REP}$  (با یادگیری از صفر) و تیم بدون یادگیری صورت دادیم که نتایج این ۱۰ بازی در شکل ۱۶ نشان داده شده است. کارایی روش یادگیری پیشنهادی (اتوماتای یادگیر  $L_{REP}$ ) در این شکل دیده می‌شود. مشاهده می‌شود که با گذشت زمان و افزایش تعداد بازیها، تیم یادگیر تدریجاً خود را با شرایط نبود ارتباطات مطمئن تطبیق می‌دهد و در جهت نیل به هدف پیش می‌رود.

در آزمایشهای دیگری ارزیابی عملکرد تیم در صورت مواجه شدن با خرابی و اشکال در عملکرد چند عامل مورد بررسی قرار گرفت. ۳ بازیکن از جناح چپ تیم اتوماتای یادگیر را برداشتیم. با توجه به آرایش ۳-۳-۴، این بازیکنان عبارت بودند از بازیکن شماره ۲ از خط دفاعی تیم (دفاع چپ)، بازیکن شماره ۶ از خط وسط تیم (هافبک چپ)، و بازیکن شماره ۱۰ از خط حمله تیم (فوروارد چپ). تیم ۸ نفری را در ۱۵ بازی متوالی آموزشی در برابر تیم بدون یادگیری که از ۱۱ بازیکن استفاده می‌کرد قرار دادیم. نتایج این بازیها در شکل ۱۷ آمده است. همانطور که نتایج نشان می‌دهد، تیم اتوماتای یادگیر ۸ نفره موفق شده است که در حین آموزش، بر نبود بازیکنان جناح چپ خود فائق آید و تدریجاً بازیها را به سود خود به پایان برساند. این امر حاکی از مناسب بودن روش یادگیری مورد استفاده می‌باشد.

در آزمایشهای دیگری، به بررسی اثر جریان باد در زمین بازی بر روی یادگیری عاملها پرداختیم. نتایج حاصله در شکل‌های ۱۷ و ۱۸ داده شده است. همانطور که مشاهده می‌شود، تیم اتوماتای یادگیر پیشنهادی موفق شده است که در حین آموزش، بر جریان باد در زمین فائق آید و تدریجاً بازیها را به سود خود به پایان برساند.

جدول ۴: میانگین آمار برای ۷ بازی آخر (پس از آموزش ۲۵ ساعته) بین تیم اتوماتای  $L_{REP}$  با تیم SALOO.

۴۵	درصد مالکیت توپ توسط تیم حریف (Saloo)
۵۵	درصد مالکیت توپ توسط تیم خودی ( $L_{REP}$ )
۲۴/۵	درصد گردش توپ در ۱/۳ زمین حریف (Saloo)
۴۲	درصد گردش توپ در ۱/۳ میانی زمین
۳۳/۵	درصد گردش توپ در ۱/۳ زمین خودی ( $L_{REP}$ )
۱۱۲/۷	ماکزیمم زمان در اختیار داشتن توپ بصورت ممتد توسط تیم حریف (Saloo)
۱۳۴/۲	ماکزیمم زمان در اختیار داشتن توپ بصورت ممتد توسط تیم خودی ( $L_{REP}$ )
۸/۸	ماکزیمم تعداد رد و بدلهای متوالی توپ توسط تیم حریف (Saloo)
۱۲	ماکزیمم تعداد رد و بدلهای متوالی توپ توسط تیم خودی ( $L_{REP}$ )

جدول ۳: میانگین آمار برای ۷ بازی اول بین تیم اتوماتای  $L_{REP}$  با تیم SALOO.

۵۲/۶	درصد مالکیت توپ توسط تیم حریف (Saloo)
۴۷/۴	درصد مالکیت توپ توسط تیم خودی ( $L_{REP}$ )
۱۰/۵	درصد گردش توپ در ۱/۳ زمین حریف (Saloo)
۴۷	درصد گردش توپ در ۱/۳ میانی زمین
۴۲/۵	درصد گردش توپ در ۱/۳ زمین خودی ( $L_{REP}$ )
۴۳۵	ماکزیمم زمان در اختیار داشتن توپ بصورت ممتد توسط تیم حریف (Saloo)
۱۱۲	ماکزیمم زمان در اختیار داشتن توپ بصورت ممتد توسط تیم خودی ( $L_{REP}$ )
۱۴	ماکزیمم تعداد رد و بدلهای متوالی توپ توسط تیم حریف (Saloo)
۸	ماکزیمم تعداد رد و بدلهای متوالی توپ توسط تیم خودی ( $L_{REP}$ )

جدول ۶: مقایسه میانگین درصد اعمال (پاس) صحیح هر بازیکن در دو روش درخت تصمیم (تیم CMUNITED98) و اتوماتای یادگیر (تیم پیشنهادی).

درصد درستی اعمال (ارزیابی پاس)	
۶۵	درخت تصمیم گیری
۷۶/۴	اتوماتای یادگیر ( $L_{REP}$ )

جدول ۵: مقایسه درصد خطای میانگین هر بازیکن تیم  $L_{REP}$  در ۷ بازی اول و ۷ بازی آخر (پس از آموزش ۲۵ ساعته) در مقابل تیم Saloo.

درصد خطا	
۴۰/۱	در ۷ بازی اول
۲۴/۶	در ۷ بازی آخر

زیادی در بالا بردن سرعت همگرایی مؤثر هستند.

### ۵-۱ آرایشهای تیمی دیگر

تا بحال، تمام شبیه‌سازیها براساس تیمهایی با آرایش ۳-۳-۴ انجام داده شد. در این قسمت به تاثیر آرایشهای تیمی دیگر بر روی همکاری عامل‌ها می‌پردازیم. بدین لحاظ غیر از آرایش ۳-۳-۴، تیمهایی مشابه ( $L_{REP}$ ) با آرایشهای ۲-۴-۴، ۱-۶-۳، ۲-۵-۳، و ۳-۴-۳ ایجاد نمودیم و یک سری بازیهای دوره‌ای بین آنها برگزار نمودیم. شکل ۲۰ اطلاعات بدست آمده از شبیه‌سازیها را نشان می‌دهد. همانگونه که نتایج نشان می‌دهد، داشتن یک شکل تیمی مناسب، در نتیجه‌گیری تیمی بسیار مؤثر است. با توجه به آمار بدست آمده، بیشترین تعداد گل زده مربوط به شکل تیمی ۲-۵-۳ و کمترین تعداد گل خورده مربوط به شکل تیمی ۲-۴-۴ است. ضمن آنکه کمترین تعداد گل زده (و یکی از بیشترین تعداد گل خورده) مربوط به شکل تیمی ۳-۳-۴ است از نتایج آزمایشها همچنین می‌توان نتیجه گرفت که در مقابل تیمهای دفاعی از یک شکل تیمی استفاده کنیم که بیشترین گل زده را دارد و یا در مقابل تیمهای حمله‌ای از شکل تیمی استفاده کنیم که آمار کمترین تعداد گل دریافتی را دارد. همچنین می‌توان نتیجه گرفت که در صورتی که نتیجه به سود تیم ماست و زمان زیادی به انتهای بازی نمانده است، شکل تیمی دفاعی (با کمترین گل خورده) را انتخاب کنیم و در صورتی که نتیجه به ضرر تیم ماست و زمان کمی تا به انتهای بازی باقی است، شکل تیمی حمله‌ای (با بیشترین گل زده) را انتخاب نماییم. در حالت عادی بازی می‌توان یک شکل تیمی که بین دو حالت فوق‌الذکر باشد استفاده کرد (مانند شکل تیمی ۳-۴-۳ که در شبه‌سازیهای این مقاله استفاده شده است).

### ۶- نتیجه گیری

تحقیق ارائه شده در این مقاله، اولین تحقیق جدی در درباره استفاده از اتوماتاهای یادگیر در همکاری در سیستم‌های چند عامله در محیط

اول و نیز ۷ بازی آخر (پس از ۲۵ ساعت آموزش) بدست آوردیم که در جدول ۵ آمده است. از نتایج جدول ۵ مشخص است که درصد انجام اعمال صحیح هر بازیکن از تیم اتوماتای یادگیر، در اثر یادگیری که در حین بازی واقع شده است بهبود یافته است.

در تیم CMUnited98 دو لایه برای رفتار چند عامله (ارزیابی پاس که بصورت غیرهمزمان و با استفاده از درخت تصمیم آموزش داده می‌شود) و رفتار تیمی (انتخاب پاس که بصورت همزمان و با استفاده از روشی بر مبنای یادگیری  $Q$  آموزش داده می‌شود) در نظر گرفته شده است [۲]. در روش پیشنهادی هر دو لایه فوق‌الذکر با هم ترکیب شده‌اند. بدین صورت که پس از تعدادی بازی آموزشی تیم را در مقابل تیم حریف در بازی آزمایشی قرار می‌دهیم، ولی چه در بازیهای آموزشی و چه در بازیهای آزمایشی روش یادگیری یکسان و همان اتوماتاهای یادگیر است و ارزیابی یک عمل از انتخاب آن عمل جدا نشده است. به همین دلیل برای مقایسه کارایی روش پیشنهادی با روش استون در تیم CMUnited98 [۲]، بایستی درصد اعمال درست و نادرست بازیکنان تیم خودی را با درصدهای پاس درست و غلط توسط عاملهای تیم CMUnited9 مقایسه کنیم. جدول ۶ درصد میانگین خطای روش پیشنهادی را با روش درخت تصمیم استون مقایسه می‌کند. توجه کنید که استون نتیجه ۶۵٪ میانگین اعمال صحیح را برای ارزیابی پاس (دومین لایه یادگیری روش لایه‌ای خود) بدست آورده است. همانگونه که این جدول نشان می‌دهد، روش پیشنهادی در این مقاله بر روش استون برتری دارد و بازیکنان پس از یادگیری کافی، درصد خطای کمتری نسبت به بازیکنان تیم CMUnited98 مرتکب می‌شوند.

نتایج بازیهای تیم اتوماتای یادگیر با تیمهای دیگر ذکر شده در ابتدای این بخش در [۱۳] آمده است. نتایج شبیه‌سازیها برای تیمهایی که از اتوماتانهای یادگیر دیگری مانند الگوریتمهای تخمین‌زن [۱۷] و یا از الگوریتم Pursuit پیوسته [۱۸] استفاده کرده‌اند نیز در [۱۳] آمده است. نتایج گزارش شده در [۱۳] نشان می‌دهند که این الگوریتمها به مقدار

[11] M. R. Meybodi and S. Lakshmirarahan, "On a class of learning algorithms which have a symmetric behavior under success and failure," *Springer-Verlag Lecture Notes in Statistics*, pp. 145-155, 1984.

[۱۲] محمدرضا خجسته و محمدرضا میبدی/ تکنیک "بهترین گوشه در مربع حالت برای عمومی سازی حالات محیطی در یک دامنه چند عامله همکاری گرا"، مجموعه مقالات هشتمین کنفرانس سالانه انجمن کامپیوتر ایران، صفحات ۴۴۶-۴۵۵، دانشگاه فردوسی مشهد، مشهد، اسفند ۱۳۸۱.

[۱۳] محمدرضا خجسته، "همکاری در سیستمهای چند عامله با استفاده از اتوماتای یادگیر"، پایان نامه کارشناسی ارشد، دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، بهار ۱۳۸۱.

[۱۴] محمدرضا خجسته و محمدرضا میبدی، "ارزیابی اتوماتای یادگیر در همکاری بین عاملها در یک سیستم چند عامله پیچیده"، مرکز تحقیقات انفورماتیک، آزمایشگاه محاسبات نرم/ دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، بهار ۱۳۸۱.

[۱۵] محمدرضا خجسته و محمدرضا میبدی، "اتوماتای یادگیر بعنوان مدلی برای همکاری در یک تیم از عاملها"، مجموعه مقالات هشتمین کنفرانس سالانه انجمن کامپیوتر ایران، صفحات ۱۲۶-۱۱۵، دانشگاه فردوسی مشهد، مشهد، اسفند ۱۳۸۱.

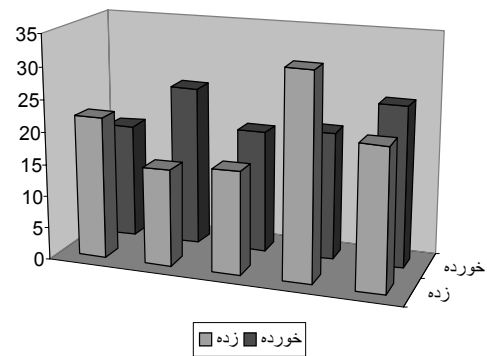
[16] I. Noda, *Team Description: Saloo*, AIST & PREST, Japan, 2001.

[17] M. A. L. Thathachar and P. S. Sastry, "A new approach to the design of reinforcement schemes for learning automata," *IEEE Trans. on Systems, Man, and Cybernetics*, vol. 15, no. 1, pp. 168-175, Jan. /Feb. 1985.

[18] B. J. Oomen and J. K. Lanctot, "Discretized pursuit learning automata," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 20, no. 4, pp. 931-938, Jul. /Aug. 1990.

محمدرضا خجسته شرح حال ایشان در زمان انتشار نشریه در دسترس نبود.

محمدرضا میبدی تحصیلات خود را در مقاطع کارشناسی و کارشناسی ارشد اقتصاد بترتیب در سالهای ۱۳۵۲ و ۱۳۵۶ از دانشگاه شهید بهشتی و در مقاطع کارشناسی ارشد و دکتری علوم کامپیوتر بترتیب در سالهای ۱۳۵۹ و ۱۳۶۲ از دانشگاه اوکلاهای آمریکا به پایان رسانده است و هم اکنون استاد دانشکده مهندسی کامپیوتر دانشگاه صنعتی امیرکبیر می باشد. نامبرده قبل از پیوستنش به دانشگاه صنعتی امیرکبیر در سالهای ۱۳۶۲ الی ۱۳۶۴ استادیار دانشگاه میشیگان غربی و در سالهای ۱۳۶۴ الی ۱۳۷۰ دانشیار دانشگاه اوهایو در ایالات متحده آمریکا بوده است. زمینه های تحقیقاتی مورد علاقه ایشان عبارتند از: الگوریتمهای موازی، پردازش موازی، محاسبات نرم و کاربردهای آن، شبکه های کامپیوتری و مهندسی نرم افزار.



شکل ۲۰: مقایسه گل های زده و خورده توسط تیمهای یادگیر اتوماتای  $L_{REP}$  با شکل های تیمی متفاوت در برابر همدیگر (از چپ به راست ۲-۴-۳، ۳-۳-۱، ۳-۵-۲، ۳-۴-۳، ۳-۵-۲).

شبیه ساز فوتبال روباتها بشمار می رود. از طریق پیاده سازی تیمهای ۲ نفره، ۵ نفره و ۱۱ نفره از عامل هایی که هر کدام از آنها مجهز به یک اتوماتای یادگیر می باشد توانایی اتوماتاهای یادگیر در ایجاد همکاری در سیستمهای چندعامله مورد بررسی و ارزیابی قرار گرفت. یک روش جدید عمومی سازی حالات محیط به نام "بهترین گوشه در مربع حالت" نیز معرفی و پیاده سازی گردید.

## مراجع

- [1] G. Weiss, *Multiagent Systems: A Modern Approach to Distributed Artificial Intelligence*, The MIT Press, London, 1999.
- [2] P. Stone, *Layered Learning in Multi-Agent Systems*, Ph.D. Thesis, School of Computer Science, Carnegie Mellon University, Dec. 1998.
- [3] I. Noda, *Team GAMMA: Agent Programming on Gaea*, in H. Kitano, editor, *RoboCup-97: Robot Soccer World Cup I*, pp. 500-507, Springer Verlag, Berlin, 1998.
- [4] RoboCup web page, at URL <http://www.robocup.org>, 1997.
- [5] H. Kitano, editor, *RoboCup-97: Robot Soccer World Cup I*, Springer Verlag, Berlin, 1998.
- [6] D. Andre *et al.*, *Soccer Server Manual*, Version 4.0, Technical Report RoboCup 1998-001, RoboCup, 1998.
- [7] K. S. Narendra and M. A. L. Thathachar, *Learning Automata: An Introduction*, Prentice-Hall Inc., 1989.
- [8] P. Mars, J. R. Chen, and R. Nambir, *Learning Algorithms: Theory and Applications*, in *Signal Processing, Control and Communications*, CRC Press, Inc., pp. 5-24, 1996.
- [9] S. Lakshmirarahan, *Learning Algorithms: Theory and Applications*, New York, Springer-Verlag, 1981.
- [10] M. R. Meybodi and S. Lakshmirarahan, " $\epsilon$ -optimality of a general class of absorbing barrier learning algorithms", *Information Sciences*, vol. 28, pp. 1-20, 1982.