

همکاری در سیستم‌های چند عامله با استفاده از اتوماتاهای یادگیر

محمد رضا خجسته و محمدرضا میدی

شناختی خود و نیز براساس تأثیر مستقیم بر ورودی‌های قابل دریافت هم‌دیگر از طریق فعالیتهای ارتباطی با هم‌دیگر همانگ 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 چند عامله می‌توان به فوتیال روباتها یا روبوکاپ^۱ اشاره کرد. محیط فوتیال روباتها 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 توب تا پیچیده‌ترین استدلال‌ها که استراتژی‌های تیم خود و تیم مقابل را روشن می‌سازند، تشکیل شوند.

به دلیل وجود پیچیدگی‌های موجود در بسترها می‌توان بستر شبیه‌سازی فوتیال روباتها و روپارویی 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^۲، نیاز به همکاری^۳ و دارای دشمن با اهداف متقابل^۴ 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 کارآیی خود را بهبود بخشنند. یکی از اهداف این مقاله بررسی کارآیی اتوماتای یادگیر در همکاری بین عاملهای عضو یک تیم در یک بستر تست شبیه‌سازی فوتیال روباتها می‌باشد. با استفاده از بستر تست شبیه‌سازی فوتیال روباتها به بررسی کارآیی اتوماتای یادگیر در همکاری بین عاملهای عضو یک تیم پرداخته شده است. با پیاده‌سازی تیمهای ۲ نفره، ۵ نفره و ۱۱ نفره از عاملهای 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، کارآیی 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 چند عامله پیچیده، داشتن روشی برای عمومی‌سازی حالات محیطی امری ضروری است انتخاب مناسب چنین روشی، در تعیین حالات و اعمال عامل نقشی تعیین کننده دارد. بهمین دلیل در این مقاله به معرفی تکنیک جدیدی بنام "تکنیک بهترین گوشه در مربع حالت" نیز می‌پردازیم. این روش فضایی عاملهایی دیگر را در خود خودش به تنها یک عمل کند ولی معمول بر آن است که یک عامل با دیگر عامل‌ها ارتباط متقابل داشته باشد. عاملهای در جهت دست یافتن به اهداف خود یا جامعه‌ای که در آن زندگی می‌کنند با یکدیگر همکاری می‌کند. وقتی که یک گروه از عاملهای در یک سیستم چند عامله در یک هدف دراز مدت سهیم باشند، آنها تشکیل یک تیم را می‌دهند.

اعضای یک تیم رفتار خود را براساس سازگار کردن فرآیندهای 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 میدی، آزمایشگاه محاسبات نرم، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، تهران، ایران (email: meybodi@ce.aut.ac.ir)

1. Learning Automaton

کلید واژه: اتوماتای یادگیر^۱، عامل، سیستم‌های چند عامله، فوتیال روباتها، همکاری.

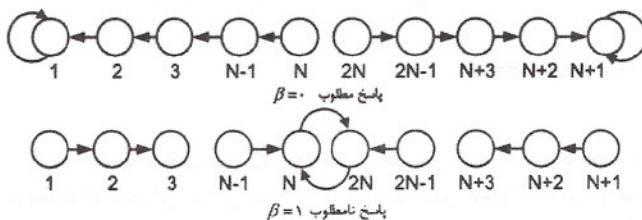
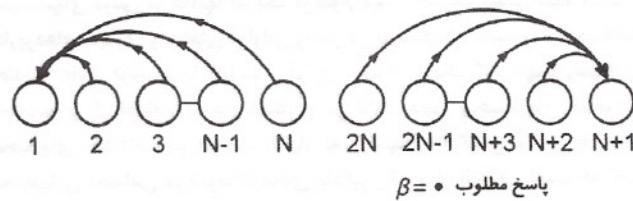
۱- مقدمه

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

اعضای یک تیم رفتار خود را براساس سازگار کردن فرآیندهای 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 میدی، آزمایشگاه محاسبات نرم، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، تهران، ایران (email: meybodi@ce.aut.ac.ir)

شکل ۲: نمودار تغییر وضعیت اتماتای $L_{2N,2}$.

شکل ۴: نمودار تغییر وضعیت اتماتای Krinsky.

ارتباط بین اتماتای یادگیر و محیط را نشان می‌دهد.

۲-۲ اتماتای یادگیر با ساختار ثابت

اتماتای یادگیر با ساختار ثابت توسط ۵ تایی $\{LA\} = \{\alpha, \beta, F, G, \phi\}$ نشان داده می‌شود که $\alpha = \{\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_r\}$ مجموعه اعمال اتماتا، $\beta = \{\beta_1, \beta_2, \beta_3, \dots, \beta_r\}$ مجموعه ورودی‌های اتماتا، $F: \phi \times \beta \rightarrow \phi$ تابعی که براساس پاسخ محیط، وضعیت جدید را می‌یابد، $G: \phi \rightarrow \alpha$ تابع خروجی که وضعیت کنونی را به خروجی بعدی می‌نگارد و $\{\phi_1, \phi_2, \dots, \phi_n\}$ مجموعه وضعیت‌های داخلی اتماتا می‌باشد.

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

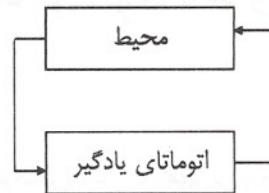
• اتماتای $L_{2N,2}$: این اتماتا تعداد پاداش‌ها و جریمه‌های دریافت شده برای هر عمل را نگهداری کرده و تنها زمانی که تعداد جریمه‌ها بیشتر از پاداش‌ها می‌گردد، عمل دیگر را انتخاب می‌کند. نمودار تغییر وضعیت این اتماتای مطابق شکل ۲ می‌باشد.

• اتماتای $G_{2N,2}$: در این اتماتا برخلاف اتماتای $L_{2N,2}$ ، عمل α_2 حداقل N بار انجام گیرد (پس از گرفتن N جریمه) تا اینکه عمل α_1 دوباره انتخاب شود. گراف تغییر وضعیت این اتماتا برای پاسخ مطلوب مانند اتماتای $L_{2N,2}$ بوده و برای پاسخ نامطلوب مطابق شکل ۳ می‌باشد.

• اتماتای Krinsky: این اتماتا زمانی که پاسخ محیط نامطلوب است، مانند $L_{2N,2}$ رفتار می‌کند. اما برای پاسخ مطلوب هر وضعیت ϕ_i ($i = 1, 2, \dots, N$) به وضعیت ϕ_{i+1} و هر وضعیت ϕ_{i+1} به وضعیت ϕ_{i+2} می‌رسد. بنابراین همیشه N پاسخ نامطلوب متواالی لازم است تا اتماتا عمل خود را عوض کند. نمودار تغییر وضعیت این اتماتا برای پاسخ نامطلوب مانند اتماتای $L_{2N,2}$ بوده و برای پاسخ مطلوب مطابق شکل ۴ می‌باشد.

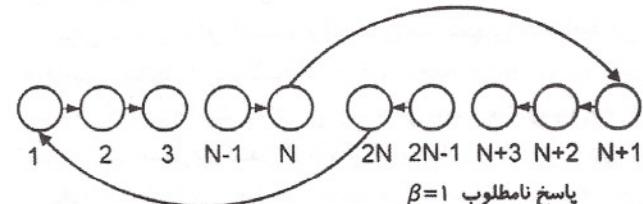
• اتماتای Krylov: در این اتماتا زمانی که پاسخ محیط مطلوب است، تغییر وضعیت مانند اتماتای $L_{2N,2}$ می‌باشد. اما زمانی که پاسخ محیط نامطلوب باشد، هر وضعیت ϕ_i ($i \neq 1, N, N+1, 2N$) با احتمال $1/5$ به وضعیت ϕ_{i+1} و با احتمال $4/5$ به وضعیت ϕ_{i-1} متنقل می‌شود.

مجموعه پاسخ مجموعه ورودی



{\alpha} مجموعه اعمال {\beta} مجموعه ورودی

شکل ۱: ارتباط بین اتماتای یادگیر و محیط.

شکل ۵: نمودار تغییر وضعیت اتماتای $G_{2N,2}$.

نظایر آن برای عمومی‌سازی فضای حالت عامل استفاده شده است.

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

۲-۲ اتماتاهای یادگیر

اتماتاهای 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 کارآبی خود را بهبود بخشنند. یک اتماتای یادگیر تعداد محدودی عمل را می‌تواند انجام دهد. هر عمل انتخاب شده توسط محیطی احتمالی ارزیابی می‌گردد و پاسخی به اتماتای یادگیر داده می‌شود. اتماتای یادگیر از این پاسخ استفاده نموده و عمل خود برای مرحله بعد انتخاب می‌کند [۷] و [۸].

اتماتاهای یادگیر به دو گروه تقسیم می‌گردند:

الف- اتماتای یادگیر با ساختار ثابت^۱

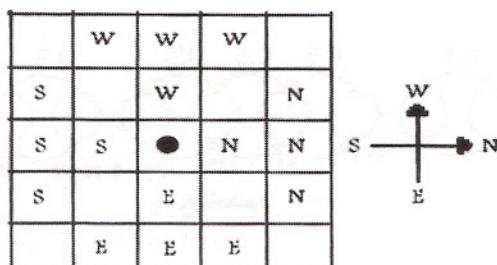
ب- اتماتای یادگیر با ساختار متغیر^۲.

۱-۲ محیط

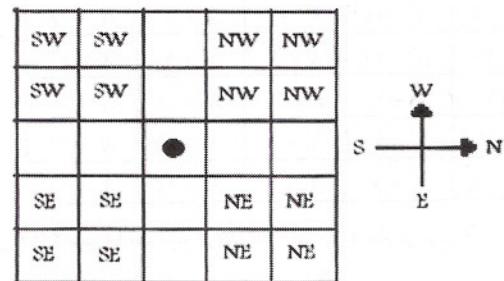
محیط را می‌توان توسط سه‌تایی $E \equiv \{\alpha, \beta, c\}$ تعریف نمود که $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_r\}$ مجموعه ورودی‌ها، $\beta = \{\beta_1, \beta_2, \dots, \beta_r\}$ مجموعه خروجی‌ها و $c = \{c_1, c_2, \dots, c_r\}$ مجموعه احتمالهای جریمه‌شدن می‌باشد.

اهرگاه β_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 α_i نتیجه نامطلوب داشته باشد می‌باشد. در محیط پایدار^۳ مقدار c_i بدون تغییر باقی می‌مانند، حال آن که در محیط ناپایدار^۴ این مقادیر در طی زمان تغییر می‌کنند. شکل ۱

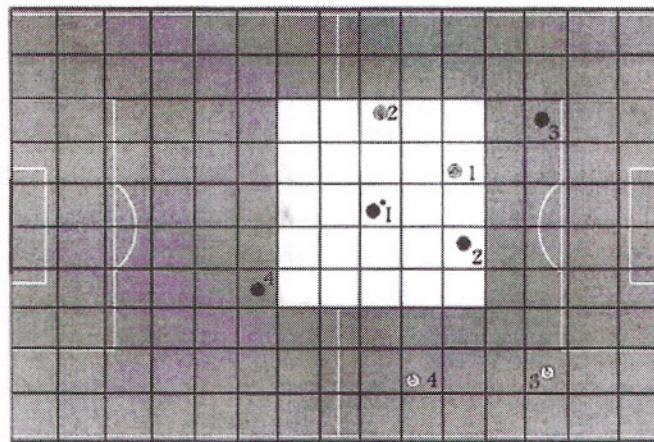
1. Fixed Structure
2. Variable Structure
3. Stationary
4. Non-Stationary



شکل ۶: 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 ۸: وضعیت بازیکنان هم تیمی و حریف با توجه به فضای محلی عامل صاحب توب در یک بازی ۴ در مقابل ۴.

می‌دهد. یعنی، عامل صاحب توب 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 فضای 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 توب، 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 توب 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 نگاشت 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 توب، 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 اتموماتیک 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 توب ببیند، موقعیت خود و دیگر بازیکنان موجود در حوزه دید خود را به بازیکن صاحب توب اعلام می‌کند تا بازیکن صاحب توب، حتی المقدور، 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-۲۱ و یا X-۲۲ و یا X-۱۲ و یا X-۱۱) ببیند، یک مقدار عددی مثبت به کمیت عددی متناظر با جهت 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 توب) در آن جهت، متناسب است. در واقع، عامل‌های نزدیک‌تر (و درون فاصله محلی تا ۲۴ مربع حالت پیرامون عامل صاحب توب) تاثیر بیشتر و عامل‌های دورتر (و درون فاصله محلی تا ۲۴ مربع حالت پیرامون عامل صاحب توب) تاثیر کمتری را بر روی ضرایب اطمینان عامل صاحب توب خواهد داشت و عامل‌های خارج از فاصله محلی (خارج از ۲۴ مربع حالت پیرامون عامل)، تاثیری در تغییر ضرایب اطمینان عامل صاحب توب 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 توب، عامل صاحب توب دارای ۸ حالت و ۸ عمل (که متناظر با ارسال توب به سمت مرکز ۸ مربع بالا فاصله پیرامون عامل می‌باشد) خواهد بود. لازم بذکر است که تقسیم زمین به مربعهای ۷×۷ متر مربعی براساس تجربه کسب شده در کار با محیط شبیه‌ساز و میزان جابجایی توب در اثر هر ضریب 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 توب در آن قرار دارد، عامل صاحب توب بر حسب آنکه عامل‌های فوق در کدام یک از ۸ جهت خود قرار دارند کمیت‌های عددی خود را تغییر

$X - 22$	$X - 12$	$X - 2$	$X + 8$	$X + 18$
$X - 21$	$X - 11$	$X - 1$	$X + 9$	$X + 19$
$X - 20$	$X - 10$	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 ۲۴: مربع اطراف مربع دوباره عامل مورد نظر در مربع X

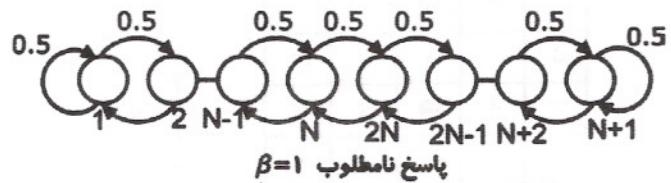
باشد، هدف پیدا کردن تابع $V \rightarrow S \rightarrow f$ می‌باشد. بدین صورت، پس از داشتن تابع f ، عامل می‌تواند از V برای یادگیری عمل مناسب خود در محیط استفاده نماید. علاوه بر این تابع، یک تابع تقسیم وظایف $P: S \rightarrow M$ نیز مورد نیاز است. تابع P مجموعه حالات دامنه را بین عاملهای موجود در محیط 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 باشد 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 فوتیال که محیط دامنه چندعامله محسوب می‌شود به 150×150 مربع یکسان با ضلع ۷ متر تقسیم می‌شود. از این طریق محیط پیوسته پیرامون عامل به محیطی گسته 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، $2^4 = 16$ مربع اطراف آن را بعنوان شاعع دید آن عامل در نظر می‌گیریم. به شکل ۶ توجه کنید.

در این شکل، مربعهای $X+10, X+11, X+1, X-1, X-9, X-10, X-11, X-12$ 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 (منتهایی چپ و بالا) تا شماره ۱۴۹ برای آخرین مربع (منتهایی راست و پایین) انجام گرفته است. هر ستون دارای $10 \times 10 = 100$ مربع و هر سطر دارای $15 \times 10 = 150$ مربع می‌باشد. در شکل‌های ۷ و ۸، مربع‌های ۸ جهت اطراف مربع دوباره عامل (با اختیار جهت شمال به سوی دروازه حریف) نشان داده 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 $X+9, X+8, X+1, X-9$ و $X+10, X+11, X+12$) بینند، یک مقدار عددی منفی به کمیت عددی متناظر با جهت شمال غرب خود اضافه می‌کند. این وضعیت برای عاملهای هم تیمی بر عکس است. به عنوان مثال، اگر عامل صاحب توب (واقع در مربع حالت X) یکی از عاملهای خودی را در درون یکی از



شکل ۵: نمودار تغییر وضعیت اتوماتای Krylov

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 اتومات، $\beta = \{\beta_1, \beta_2, \dots, \beta_m\}$ مجموعه ورودی‌های اتومات، $p = \{p_1, p_2, \dots, p_r\}$ 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 اتومات‌ها، اگر عمل i در مرحله n انتخاب شود و پاسخ مطلوب از محیط دریافت نماید، احتمال $p_i(n)$ افزایش یافته و سایر احتمالها کاهش می‌یابند و برای پاسخ نامطلوب احتمال $p_i(n)$ کاهش یافته و سایر احتمالها 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 a از b بزرگتر باشد، الگوریتم را L_{REP} می‌نامیم و هنگامی که b بزرگتر 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 اتومات 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 فوتیال 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

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

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

۶	۵	۴	۳	۲	۱	
۵-۱۶	۵-۱۲	۴-۹	۴-۷	۳-۳	۲-۱	$L_{2N,2}$
۵-۱۶	۳-۱۳	۲-۱۰	۲-۸	۱-۵	۱-۲	$G_{2N,2}$
۴-۱۷	۳-۱۵	۳-۱۱	۳-۷	۲-۴	۱-۲	Krinsky
۳-۱۲	۳-۱۱	۳-۹	۲-۶	۱-۴	۱-۲	Krylov
۴-۱۷	۳-۱۵	۳-۱۳	۲-۱۰	۲-۵	۱-۳	Q

همکاری مورد ارزیابی قرار دهیم.

۴-۱ شبیه‌سازیها برای تیم ۲ نفره

چندین سری شبیه‌سازی انجام گرفته است. در اولین سری شبیه‌سازیها، به پیاده‌سازی تیمهای ۲ نفره از عاملها می‌پردازیم. این سری از شبیه‌سازیها با دو روش برای تعیین حالت هر عامل در محیط خود انجام می‌گیرد. یکی از این دو روش یک عمومی‌سازی ساده و دیگر روش، تکیک "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 توب و ۴ حالت برای بازیکن بدون توب خلاصه می‌شود [۱۳]. برای هریک از این چهار حالت یک اتماتای یادگیر با ساختار ثابت و با عمق حافظه ۳ درنظر گرفته می‌شود. هر اتماتا، 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 اتماتای یادگیری در مقایسه با یک تیم بدون یادگیری به سرعت یاد می‌گیرد که در چه حالاتی، باستی 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 اتماتای یادگیر می‌تواند در حین بازی خود را با نحوه بازی تیم حریف تا حد زیادی تطبیق دهد. تعداد ۵۰ بازی را بین تیمهای اتماتای 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 بدون یادگیری عدد سمت چپ نشان‌دهنده تعداد گلهای زده توسط تیم بدون یادگیری می‌باشد. همانگونه که مشاهده می‌شود کلیه بازیها به سود تیمهای یادگیری 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. همان‌گونه که این جدول نشان می‌دهد، تفاوت زیادی بین اتماتاهای یادگیر مختلف در سری شبیه‌سازیها انجام شده مشاهده نمی‌شود، هرچند که روش‌های یادگیری Q و اتماتای Krinsky اندکی بهتر از دیگر اتماتاهای عمل کرده است.

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

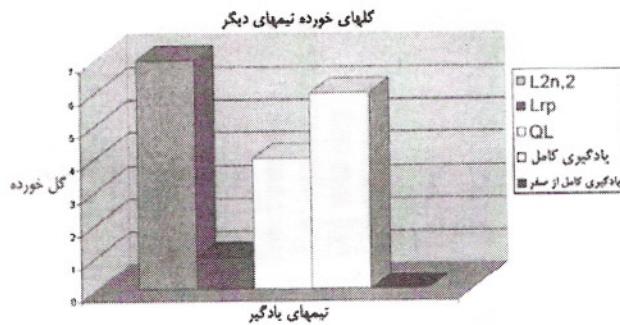
محدود است و هر بازیکن کنار و پشت خود را نمی‌بیند، برای بازیکن صاحب توب جیاتی است و لذا در این مورد (و تنها در این مورد) از امکان شناوری (و آنهم بصورت محدود) استفاده شده است. در واقع در چنین حالتی، بازیکن بدون توب، مدل دنیای خود را به بازیکن صاحب توب هم تعیین خود اعلام می‌کند و وی را در تعیین حالت و به تبع آن، تعیین عمل مناسب در آن حالت، یاری می‌کند. بازیکن بدون توب نیز، به مرکز مربع بلاfaciale‌ی (از مربع حالت در بردارانده خود) می‌رود که به منظور دریافت توب (در صورت لزوم) از بازیکن هم تعیین صاحب توب، مناسب تشخیص می‌دهد. بدین نحو، در شبیه‌سازیها انجام گرفته، بازیکن بدون توب نیز می‌تواند دارای حالات و اعمال خاص (حرکت به سمت یکی از ۸ جهت اطراف خود به منظور دریافت توب) باشد. یعنی بازیکن می‌تواند "حرکت بدون توب" انجام دهد. باید توجه داشت که عمل مناسب انتخابی عامل صاحب توب در هر حالت، الزاماً به معنی پاس و یا شوت نمی‌باشد و در بعضی موارد، ممکن است دریبل و یا نگهداشتن موقع توب به منظور یافتن روزنای 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 توب، پاس به عقب، و حتی بازکردن بازی و ارسال توب به فضای خالی (با احتمال دریافت مناسب توب برای بازیکن هم تعیین و یا خود در چند سیکل بعد) نیز استفاده می‌کند [۱۳].

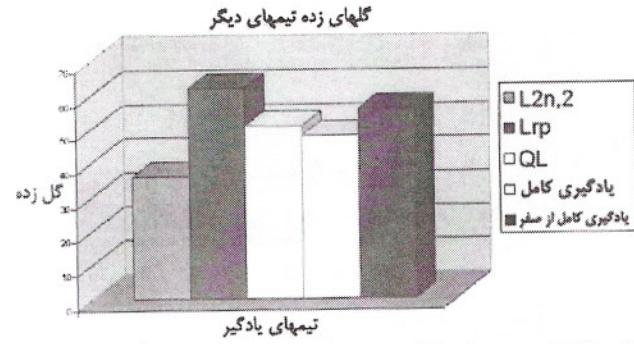
لازم به ذکر است که یادگیری در کلیه شبیه‌سازیها انجام گرفته کاملاً چند عامله و توزیع شده است و برخلاف روش‌های موجود، بازیکن صاحب توب (بدون توب)، قصد خود را از ارسال (دربافت) توب به (از) بازیکن بدون توب (صاحب توب) اعلام نمی‌دارد و بدین ترتیب، هر عامل در انتخاب عمل بهینه خود کاملاً خودمختار است و درین انجام بازی بدون ایجاد ارتباط با دیگر بازیکنان سعی می‌کند که با آنها همکاری کنند. شبیه‌سازیها 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 اتماتای یادگیر

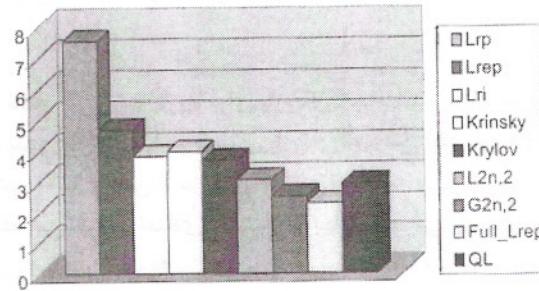
تا بحال از روش‌های مختلفی از جمله یادگیری Q، الگوریتمهای ژنتیک، درختهای تصمیم‌گیری و یادگیریهای رفتاری [۲] برای یادگیری عامل‌های فوتbalیست استفاده شده است. در این بخش از مقاله، توانایی اتماتای یادگیر بمنظور ایجاد همکاری بین عامل‌های موجود در یک تیم روبوکاپ برای رسیدن به یک هدف مشخص تعیین که همانا برداش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



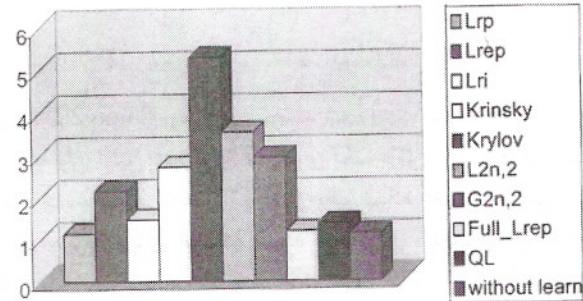
شکل ۱۱: گلهای خورده تیمهای یادگیر ۵ نفره در ۱۰ شبیه سازی (در حین آموزش).



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 پادگیری با تیمهای دیگر.



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

نایابیهای ۲ نفره در هیچ دیداری بازنشسته نبودند. نتایج آزمایشها در اشکال ۱۱ آمده است.

۳-۴- شبیه‌سازی‌ها برای تیم ۱۱ نفره

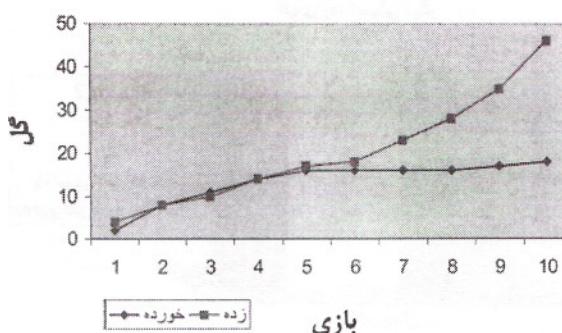
شبیه سازی های ارایه شده در این قسمت به بررسی یادگیری در نحوه همکاری در بین عامل های موجود در یک تیم ۱۱ نفره و مقایسه ان با دیگر نیمه ها می پردازد. با توجه به این که در بازی های ۱۱ نفر در مقابل ۱۱ نفر، ممکن است تعداد گلهای رد و بدل شده خیلی زیاد نباشد به منظور بررسی نتایج شبیه سازی ها در این سری (و سری های بعد)، از معیار های دیگری هم بمنظور نشان دادن کارایی استفاده می شود. پس از بررسی معیار های مطرح در همکاری بین بازی کنان یک تیم فوتبال، معیار های زیر نظر مناسب می باشد.

- درصد مالکیت توب توسط تیم خودی در مقایسه با مورد مشابه در نیم حریف در حین بازی.
 - درصد گردش توب در $\frac{1}{3}$ زمین خودی، $\frac{1}{3}$ میانی زمین، و $\frac{1}{3}$ زمین حریف در حین بازی.
 - حداکثر زمان در اختیار داشتن توب بصورت ممتد توسط تیم که بر حسب سیکا سنجیده م شهد.

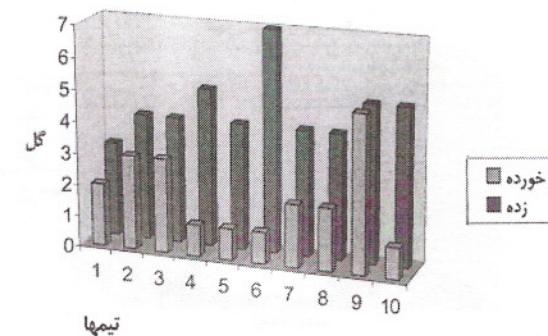
- حداقل تعداد رد و بدلهای متواالی توب بدون برخورد با حریف.
- میانگین درصد خطای (اعمال) بازیکنان تیم خودی در حین بازی.
- برای سازماندهی ۱۱ بازیکن درون زمین برای هر تیم آرایش تیمی ۴-۳-۳ انتخاب گردیده است. همانند شبیه‌سازیهای قبلی، در این سری از آزمایشها نیز برای مقایسه روند یادگیری در تیمهای یادگیر، یک 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 اتوماتی 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 توب بعد از خاتمه این عمل هنوز تیم اوست و توب به سود تیم او جلوتر رفته است (بطور مثال تیم او گل زده است) و در غیر این صورت بازیکن جرمیه می‌شود. همچنین اگر انجام دهنده آخرين عمل، خود من بوده‌ام (یا اگر هنوز نتیجه آخرين عمل خود را ندیده‌ام) و تیم صاحب توب، تیم حریف است و اگر فاصله موقعیت قبلی توب تا موقعیت فعلی توب خیلی زیاد نیست (بطور مثال به تیم حریف گل نزنده‌ام) بازیکن جرمیه می‌شود و در غیر این صورت به او پاداش داده می‌شود. نتیجه عمل که پاداش یا جرمیه 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

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 نفرات 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 روبوکاپ با روباتهایی با اندازه‌های متوسط و بزرگ و نیز فوتbal داخل سالنی بوده است. در شبیه‌سازیهای ۵ نفره، تیمهای یادگیر یادگیر با ساختار ثابت، L_{RN} ،^۲ بعنوان نماینده تیمهای یادگیر با ساختار متغیر، Q بعنوان L_{RP} بعنوان نماینده تیمهای دارای روش یادگیری به غیر از اتموماتیک یادگیر، و نماینده‌ای از تیمهای دارای استعدادهای متفاوت به غیر از اتموماتیک 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 از زمین که در واقع محدوده‌ای از مربعهای حالت می‌باشد مشخصی از زمین که در بازیهای انجام شده، تیمهای یادگیر، همانند تخصیص داده شده است. در بازیهای انجام شده، تیمهای یادگیر، همانند



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



شکل ۱۶: مقایسه گلهای خورده (میله های جلو) و گلهای زده (میله های عقب) تیمهای مختلف در مقابل تیم دستتویس (ثابت)، از چپ به راست: بدون یادگیری، Q، کامل $G_{2N,2}$ ، Krinsky، Krylov، L_{RN} ، L_{RP} ، L_{RI} ، L_{REP} ، L_{REP} بازیها را در مقابل تیم "نگاشت ثابت" خلاصه می کند.

تیم 2001 Saloo توسط نودا نوشته شده است. نودا در شرح تیم خود تقلید را عنوان اولین قدم برای سازگاری یک عامل با عامل های دیگر در محیطی چند عامله می دارد [۱۶]. او از یک شبکه عصبی بازگشتی برای یادگیری و برنامه ریزی استفاده نموده است. این شبکه عصبی دارای ۲ وظیفه است: پیشگویی محیط و شناخت بازی. شبکه عصبی استفاده شده می تواند عنوان یک تمیزدهنده بین انواع بازیها آموزش داده شده مورد استفاده قرار گیرد. وی نتیجه گیری کرده است که این شبکه به همراه معماری وی می توانند انواع گوناگونی از بازیها را مجدد تولید کند [۱۶].

در آزمایشها این قسمت از تیم اتوماتای L_{REP} و تیم یادگیر کامل L_{REP} استفاده شده است. این انتخاب بدین دلیل بوده است که این تیمهای در شبیه سازیها قبلی [۱۲] تا [۱۵] نتایج خوبی را تولید کرده است. در اولین سری از شبیه سازیها به ارزیابی تیم اتوماتای یادگیر در شرایطی که نویز در محیط وجود داشته باشد می پردازیم. پارامترهای متعددی در کارگزار شبیه سازی فوتbal وجود دارند که می توان 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 شبیه سازیها، اثر پارامترهای rand که معرف مقدار نویز در موارد مختلف محدود بررسی قرار داده شده است. بدینصورت که پارامتر player_rand را از مقدار $0.1 / 0.2$ ، پارامتر ball_rand را از $0.05 / 0.1$ ، و بالاخره پارامتر kick_rand را از $0.1 / 0.05$ تغییر دادیم. پارامتر اول باعث ایجاد اختلال (نویز) توسط کارگزار در حرکت بازیکن می شود و پارامترهای دوم و سوم، همین نقص را در مورود حرکت توپ در زمین و زدن ضربه به توپ بر عهده دارند. ۱۰ بازی متواالی بین تیم L_{REP} (با یادگیری از صفر) و تیم بدون یادگیری برگزار نمودیم. هدف از این شبیه سازی این بود که بینیم آیا روش های یادگیری پیشنهادی در این مقاله می توانند در این گونه شرایط نیز موثر واقع شوند. شکل ۱۵ نتایج را (تصویر تجمعی) نشان می دهد. در این شکل، کارآیی روش یادگیری پیشنهادی (اتوماتای یادگیر L_{REP}) به صورت روشن مشاهده می شود. همانگونه که شکل نشان می دهد، با افزایش تعداد بازیها، تیم اتوماتای یادگیر موفق می شود که رفتہ رفتہ 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 می تواند خود را با شرایط محیطی سازگار کند. گرچه بدليل تاکید بر خصوصیت خود مختاری عامل های تیمهای پیشنهادی از حداقل ارتباطات بین بازیکنان در حین بازی سود می برد ولی بدليل اینکه کاتال ارتباطی بین بازیکنان شلوغ و فاقد اطمینان است، برای همین حجم اندک ارسال (و نه تبادل) اطلاعات اتکا به ارتباطات بین عامل ها در چنین دامنه های قابل اطمینان نیست. تیمهای اتوماتاهای

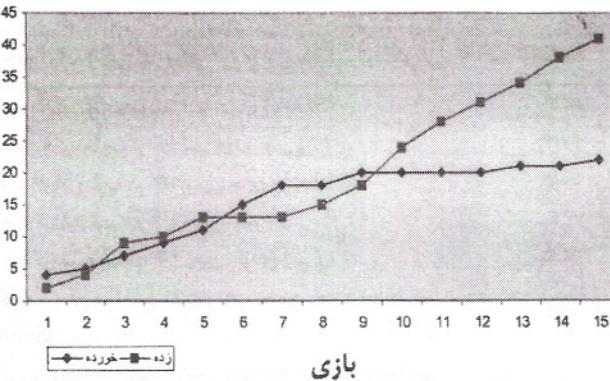
و ۱۳ نتایج کلی شبیه سازیها برای بازیهای انجام شده بین تیمهای یادگیر و تیم بدون یادگیری را در دوسری بازیهای آموزشی (۱۵ بازی) و بازیهای آزمایشی (۳ بازی پس از بازیهای آموزشی) خلاصه می کنند. همانگونه که مشاهده می شود با تعدادی محدود بازیهای آموزشی تیمهای اتوماتای یادگیر با ساختار ثابت به برتری نسبی می رسد. در مورد بازیهای آزمایشی (پس از آموزش و تنها با استخراج مقادیر یادگرفته شده) به نظر می رسد که برتری نسبی (در غلبه بر تیم بدون یادگیری) با تیمهای اتوماتای یادگیر با ساختار متغیر می باشد. جزییات این آزمایشها در [۱۳] آمده است. شکل ۱۴ گلهای زده و خورده تیمهای یادگیر را در مقابل تیم "نگاشت ثابت" خلاصه می کند.

۵- تستهای ارزیابی برای تیمهای مبتنی بر اتوماتای یادگیر

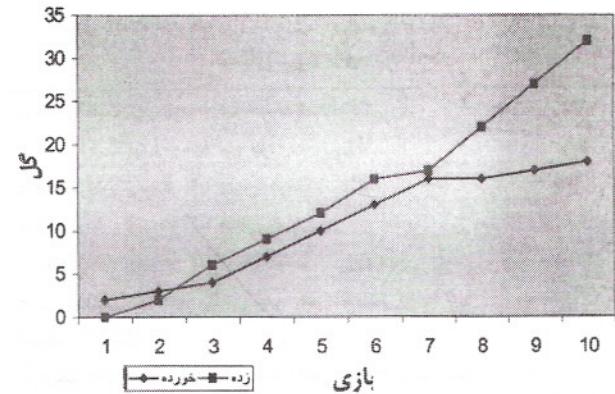
در بخشها قبلي، عملکرد اتوماتاهای یادگیر را در انجام یک کار گروهی در بین عاملهای عضو یک تیم مورد بررسی قرار دادیم. تمام شبیه سازیهای انجام شده در موارد فوق، در شرایط طبیعی دامنه برگزار شدند. در این بخش به ارزیابی تیمهای یادگیر از طریق بازی با بعضی از تیمهای شرکت کننده در مسابقات روبوکاپ جهانی می پردازیم.

گرچه تیم اتوماتای یادگیر بدليل پیاده سازی نکردن مواردی مانند استراتژی حمله، استراتژی دفاع، مدل کردن ... تغییر استراتژی بازی در زمانهای خاص، مرتب و ... یک تیم مسابقه محسوب نمی شود ولی بازی با تیمهای شناخته شده که در مسابقات روبوکاپ شرکت می کنند، می تواند ما را در ارزیابی مدل اتوماتاهای یادگیر در یادگیری عامل ها یاری کند. لازم به ذکر است که تیم اتوماتای یادگیر یک تیم تحقیقاتی می باشد و صرفاً برای مطالعه و بررسی توانایهای اتوماتاهای یادگیر در همکاری بین عامل ها در انجام یک کار گروهی استفاده شده است.

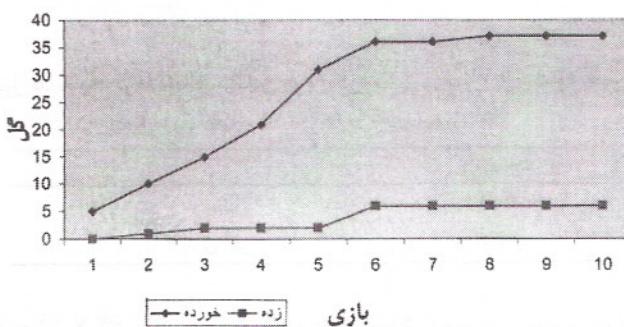
یادآوری می شود که کد پایه که برای پیاده سازی تیمهای اتواتهای یادگیر مورد استفاده قرار گرفته است، کد CMUnited98 [۲] می باشد. لازم به ذکر است که این تیم در مقایسه با تیمهای جدیدتر از مهارت های Sharif Arvand، FuzzyFoo 2001، 2000، Yberoos 2000، 2001 و Saloo 2001 که از نظر فردی اختلاف زیادی با تیم ما نداشتند (هر چند بهتر از تیم ما بودند) می توانند عنوان حریف بازیها انتخاب شوند. از بین چهار تیم فوق، تیم Saloo 2001 بدليل نزدیکتر بودن مهارت های فردی با تیم خودی مناسب ترین تیم برای مسابقه به منظور ارزیابی روش اتوماتای یادگیر تشخیص داده شد و بهمین دلیل بازیها با این تیم انجام گرفت [۱۳]. برای مشاهده نتایج بازیها با دیگر تیمهای فوق می توان به [۱۳] مراجعه کرد.



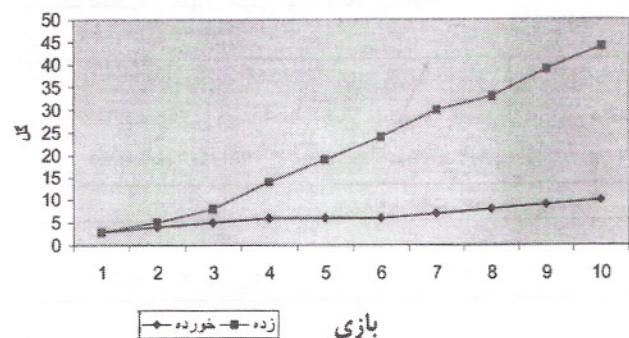
شکل ۱۷: مقایسه گلهای زده و خورده در طی ۱۰ بازی آموزشی متوالی با برداشتن ۳ بازیکن جناح چپ از تیم یادگیر زمانیکه تیم اتوماتای L_{REP} با تیم بدون یادگیری بازی می‌کند.



شکل ۱۸: 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 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 ۲۰: 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 تفاصل گل نسبتاً بالای (میانگین گل زده ۵/۷ گل در مقابل گل خورده ۰/۰ 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 مطلقی 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 بازی L_{REP} که نتایج ۷ بازی اول آنها در جدول ۳ آمده است)، تیم اتوماتای یادگیر 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 ۱/۰ 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.

یک دیگر از معیارهایی که برای ارزیابی کارآیی همکاری در بین عاملهای تعريف 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. کارآیی 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

در آزمایشها دیگری ارزیابی عملکرد تیم در صورت مواجه شدن با خرابی و اشکال در عملکرد چند عامله مورد بررسی قرار گرفت. ۳ بازیکن از جناح چپ تیم اتوماتای یادگیر را برداشتیم. با توجه به آرایش ۴-۳-۳، 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 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 آزمایشها 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 توب توسط تیم حریف (Saloo)
۵۵	درصد مالکیت توب توسط تیم خودی (L_{REP})
۲۴/۵	درصد گردش توب در ۱/۳ زمین حریف (Saloo)
۴۲	درصد گردش توب در ۱/۳ میانی زمین
۳۳/۵	درصد گردش توب در ۱/۳ زمین خودی (L_{REP})
۱۱۲/۷	ماکریم زمان در اختیار داشتن توب بصورت ممتد توسط تیم حریف (Saloo)
۱۳۴/۲	ماکریم زمان در اختیار داشتن توب بصورت ممتد توسط تیم خودی (L_{REP})
۸/۸	ماکریم تعداد رد و بدلهای متواالی توب توسط تیم حریف (Saloo)
۱۲	ماکریم تعداد رد و بدلهای متواالی توب توسط تیم خودی (L_{REP})

جدول ۵: مقایسه میانگین درصد اعمال (پاس) صحیح هر بازیکن تیم L_{REP} در ۷ بازی اول و ۷ بازی آخر (پس از آموزش ۲۵ ساعته) در مقابل تیم SALOOجدول ۳: میانگین آمار برای ۷ بازی اول بین تیم اتماتاتی L_{REP} با تیم SALOO

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

جدول ۶: 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 خطأ	در ۷ بازی اول	در ۷ بازی آخر
۴۰/۱	در ۷ بازی اول	در ۷ بازی آخر
۲۴/۶		

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 اتماتاتی یادگیر (تیم پیشنهادی)، بدین بازی واقع شده است بهمود یافته است.

در تیم CMUnited98 دو لایه برای رفتار چند عامله (ازبیابی پاس که بصورت غیرهمزان 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 \mathcal{Q} آموزش داده می‌شود) در نظر گرفته شده است [۲]. در روش پیشنهادی هر دو لایه فوق الذکر 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 اتماتاتی یادگیری است و ازبیابی یک عمل از انتخاب آن عمل جدا نشده است. به همین دلیل برای مقایسه کارآیی 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 روش 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 ازبیابی 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 اتماتاتی یادگیر با تیمهای دیگر کشیده در ابتدای این بخش در [۱۳] آمده است. نتایج شبیه‌سازیها برای تیمهایی که از اتماتاتهای یادگیر دیگری مانند الگوریتمهای تخمین زن [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

زیادی در بالا بردن سرعت همگرایی مؤثر هستند.

۱-۵ آرایشهای تیمی دیگر

تا بحال، تمام شبیه‌سازیها براساس تیمهایی با آرایش ۴-۳-۳ انجام داده شد. در این قسمت به تاثیر آرایشهای تیمی دیگر بر روی همکاری عامل‌ها می‌پردازیم. بدین لحاظ غیر از آرایش ۴-۳-۳، تیمهایی مشابه (یعنی L_{REP}) با آرایشهای ۲-۴-۲، ۱-۶-۳، ۲-۵-۳، 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 موثر است. با توجه به آمار بدست آمده، بیشترین تعداد گل زده مربوط به شکل تیمی ۴-۳-۵-۲ و کمترین تعداد گل خورده مربوط به شکل تیمی ۴-۴-۲ است. ضمن آنکه کمترین تعداد گل زده (و یکی از بیشترین تعداد گل خورده) مربوط به شکل تیمی ۴-۳-۳ 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 ضرر 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 ۳-۴-۳ که بین دو حالت فوق‌الذکر باشد استفاده کرد (مانند شکل تیمی ۴-۳-۳-۳ که در شبیه‌سازیها این مقاله استفاده شده است).

۶- نتیجه گیری

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

[11] M. R. Meybodi and S. Lakshminarahan, "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 میدی / تکنیک "بهترین گوشه در مربع حالت برای عمومی‌سازی حالات محیطی در یک دامنه چند عامله همکاری گرا"، مجموعه مقالات هشتمین کنفرانس سالانه انجمان کامپیوتر ایران، صفحات ۴۵۵-۴۴۶، 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 میدی، "ازیابی 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 میدی، "اتوماتای یادگیر بعنوان مدلی برای همکاری در یک تیم از عاملها"، مجموعه مقالات هشتمین کنفرانس سالانه انجمان کامپیوتر ایران، صفحات ۱۲۶-۱۱۵، دانشگاه فردوسی مشهد، مشهد، اسفند ۱۳۸۱.

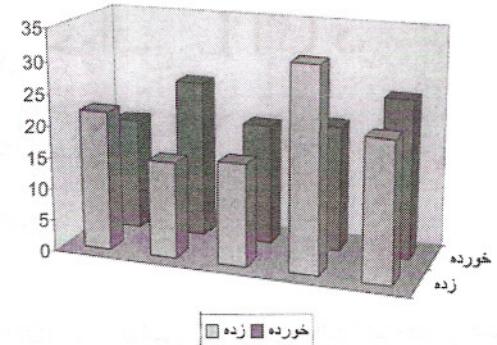
[16] I. Noda, *Team Description: Saloo*, AIST & PREST, Japan, 2001.

[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 میدی 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 اوکلاهوما امریکا به پایان رسانده است و هم اکنون استاد دانشکده مهندسی کامپیوتر دانشگاه صنعتی امیرکبیر می‌باشد. نامبرده قبل از پیوستن به دانشگاه صنعتی امیرکبیر در سالهای ۱۳۶۲ الی ۱۳۶۴ استادیار دانشگاه میسیگان غربی و در سالهای ۱۳۶۴ الی ۱۳۷۰ دانشیار دانشگاه اوهاوی در ایالات متحده امریکا بوده است. زمینه‌های تحقیقاتی مورد علاقه ایشان عبارتند از: الگوریتمهای موازی، پردازش موازی، محاسبات نرم و کاربردهای آن، شبکه‌های کامپیوتری و مهندسی نرم افزار.



شکل ۲۰: مقایسه گل‌های زده و خورده توسط تیمهای یادگیر اتوماتای L_{REP} با شکلهای تیمی متفاوت در برابر همدیگر (از چپ به راست ۳-۴-۳، ۴-۳-۲، ۳-۶-۱، ۳-۴-۳، ۳-۵-۲).

شبیه‌ساز فوتبال روباتها بشمار می‌رود. از طریق پیاده‌سازی تیمهای ۲ نفره، ۵ نفره و ۱۱ نفره از عامل‌هایی که هر کدام از آنها مجهز به یک اتوماتای یادگیر می‌باشد توانایی اتوماتاهای یادگیر در ایجاد همکاری در چند عامله چند عامله مورد بررسی و ارزیابی قرار گرفت. یک روش جدید عمومی‌سازی حالات محیط به نام "بهترین گوشه در مربع حالت" نیز معرفی و پیاده‌سازی 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. Lakshminarahan, *Learning Algorithms: Theory and Applications*, New York, Springer-Verlag, 1981.
- [10] M. R. Meybodi and S. Lakshminarahan, " ε -optimality of a general class of absorbing barrier learning algorithms", *Information Sciences*, vol. 28, pp. 1-20, 1982.