

# بهبود عملکرد حمله در تیم ربات‌های فوتبالیست با استفاده از یادگیری تقویتی

مینا خاکسار<sup>۱</sup>، دانشجو؛ ولی درهمی<sup>۲</sup>، دانشیار؛ مهدی رضائیان<sup>۳</sup>، استادیار

۱- گروه مهندسی کامپیوتر - پردیس فنی و مهندسی - دانشگاه یزد - یزد - ایران - khaksarmina@stu.yazd.ac.ir

۲- گروه مهندسی کامپیوتر - پردیس فنی و مهندسی - دانشگاه یزد - یزد - ایران - vderhami@yazd.ac.ir

۳- گروه مهندسی کامپیوتر - پردیس فنی و مهندسی - دانشگاه یزد - یزد - ایران - mrezaeian@yazd.ac.ir

**چکیده:** به دلیل عدم امکان پیش‌بینی همه وضعیت‌های ممکن برای عامل‌ها در یک سیستم چندعامله‌ی پویا و گسترده، روش‌های یادگیری ماشین، ابزار مناسبی برای کنترل رفتار عامل‌ها می‌باشد. فوتبال شبیه‌سازی شده ربات‌ها یک مسئله شناخته‌شده برای ارزیابی الگوریتم‌های یادگیری ماشین روی سیستم‌های چندعامله است. در این مقاله الگوریتم یادگیری کیو-وی (یکی از الگوریتم‌های معروف یادگیری تقویتی) جهت بهبود عملکرد حمله در تیم ربات‌های فوتبالیست دو بعدی بکار گرفته شده است. سیگنال تقویتی براساس اینکه بازیکنانی که در حمله دخالت دارند، توپ را جلوی دروازه می‌رسانند، یا اینکه توپ را از دست می‌دهند، تعریف شده است و به ترتیب عامل‌ها با توجه به وضعیت ذکر شده، جایزه و جریمه دریافت می‌کنند. جهت بهبود عملکرد از ایده تقسیم سیگنال تقویتی متناسب با مقدار خبرگی عامل‌ها در یک سیستم چندعامله استفاده شده است. در اینجا میزان خبرگی متناسب با تفاوت ارزش عمل با بالاترین ارزش با ارزش عمل با کمترین مقدار ارزش در هر حالت تعریف می‌شود. نتایج شبیه‌سازی نشان می‌دهد، بهره‌گیری از ایده تقسیم سیگنال تقویتی براساس خبرگی در آموزش باعث افزایش سرعت در آموزش و بهبود عملکرد عامل‌ها شده است. واژه‌های کلیدی: فوتبال شبیه‌سازی شده ربات‌ها، یادگیری تقویتی، یادگیری کیو-وی، سیستم‌های چندعامله، حمله.

## Improve performance of attack in the team robots soccer using reinforcement learning

Mina Khaksar<sup>1</sup>, Student; Vali Derhami<sup>2</sup>, Associate Professor; Mehdi Rezaeian<sup>3</sup>, Assistant Professor

1- Computer Engineering Department, Faculty of Engineering, Yazd University, Yazd, Iran, Email: khaksarmina@stu.yazd.ac.ir

2- Computer Engineering Department, Faculty of Engineering, Yazd University, Yazd, Iran, Email: vderhami@yazd.ac.ir

3- Computer Engineering Department, Faculty of Engineering, Yazd University, Yazd, Iran, Email: mrezaeian@yazd.ac.ir

**Abstract:** Due to the impossibility of predicting all possible states for agents in a wide dynamic multi-agent system, machine learning methods are useful tools to control agent behavior. Simulated Robot Soccer is a well known multi agent benchmark to evaluate machine learning algorithms. In this paper, QV-Learning algorithm (a well known reinforcement learning algorithm) is used to improve the performance of the attack in 2D robots soccer team. The reinforcement signal is defined based on the players involved in the attack can reach the ball in front of goal or lose the ball; They receive positive and negative reward according to the mentioned status, respectively. We use the idea of division the reinforcement signal proportional to the amount of expertness (knowledge) of agents to improve the performance. Here, the expertise is defined as the difference between highest action value and lowest action value in the each state. The simulation results show using the idea of expertise improves the train speed and the performance.

**Keywords:** Simulated Robot Soccer, reinforcement Learning, QV-Learning, multi-Agent System, attack.

تاریخ ارسال مقاله: ۱۳۹۵/۰۹/۲۹

تاریخ اصلاح مقاله: ۱۳۹۶/۰۱/۰۸ و ۱۳۹۶/۰۶/۰۶

تاریخ پذیرش مقاله: ۱۳۹۶/۰۸/۰۱

نام نویسنده مسئول: ولی درهمی

نشانی نویسنده مسئول: ایران - یزد - صفائیه - بلوار دانشگاه - دانشگاه یزد - ساختمان پردیس فنی و مهندسی ۱.

## ۱- مقدمه

از دیدگاه هوش مصنوعی توزیع شده، مسابقه فوتبال ربات‌ها نمونه‌ای از مسائل چندعامله است که مباحث تحقیقاتی متعدد را می‌توان در آن یافت [۱]. هر تیم مجموعه‌ای از ربات‌ها است که باید به‌منظور رسیدن به هدف نهایی و مشترک که همان پیروزی در بازی است با یکدیگر همکاری داشته باشند. به‌ثمر رساندن گل و جلوگیری از گل‌زدن توسط حریف را می‌توان به‌عنوان زیرمجموعه‌ای از هدف اصلی پیروزی در نظر گرفت. عامل‌ها در شرایطی که درک درستی از محیط ندارند و دید محلی نسبت به اطراف دارند، باید برای بهینه‌کردن میزان کارایی خود در محیط تلاش کنند. به‌علاوه عامل‌ها باید در یک محیط پویا به‌صورت بی‌درنگ و خصمانه عمل کنند. در نتیجه موفقیت تیمی وابسته به سرعت عملکرد و انعطاف‌پذیری عامل‌ها است. از آنجایی که فضای حالات ممکن برای بازی بسیار بزرگ است، استفاده از یک روش یکنواخت ممکن نیست و لذا طراحی استراتژی‌های مناسب در دستیابی به موفقیت ضروری می‌باشد.

توسعه روش‌های مؤثر برای هماهنگی سیستم‌های چندعامله در محیط‌های خصمانه یکی از جذاب‌ترین چالش‌های علمی است که توسط ربات‌های ترویج داده شده‌است و به‌طور کلی توسط لیگ‌های شبیه‌ساز فوتبال حمایت می‌شود. هدف کلی از مکانیسم هماهنگی در این لیگ‌ها کنترل مناسب یک تیم از بازیکنان و داشتن اختیار برای بردن مسابقه در مقابل تیم حریف است.

در محیط چندعامله فوتبال، تصمیم‌گیری برای مهارت‌های سطح بالا توسط عامل‌ها بر عملکرد تیمی تأثیرگذار است. اگرچه هنوز واضح نیست که بیشتر کدام تکنولوژی پایه مانند برنامه‌ریزی تئوری تصمیم‌گیری، یا مهارت‌های سطح پایین، مهم‌ترین نقش را در موفقیت یک تیم ربات‌ها دارد. از طرفی استراتژی‌های موفق و تاکتیک‌های فوتبال تأثیر بسیاری بر مهارت‌های بازیکن دارد. برای اینکه عامل‌ها یک استراتژی را با موفقیت عملی نمایند، باید توانایی پاس‌کاری به تمام نقاط، موقعیت‌یابی مناسب برای حمله<sup>۲</sup> و دفاع<sup>۳</sup>، علامت‌گذاری<sup>۴</sup> حریف و مهارت‌های سطح پایین‌تر را دارا باشند. از جمله مهارت‌های سطح بالای چندعامله در فوتبال ربات‌ها استراتژی، تاکتیک و شکل‌گیری، مهارت تهاجمی و مهارت تدافعی است. از آنجایی که این مقاله روی مهارت تهاجمی بازیکنان تمرکز دارد، توضیح مختصری از آن در زیر آمده‌است.

حمله به معنای گرفتن توپ و حرکت با توپ به سمت جلو و به سمت دروازه حریف است. استراتژی حمله معمولاً به‌این‌صورت تعریف می‌شود که گردش توپ بین هم‌تیمی‌ها تا زمان فعال‌شدن تابع شوت‌کردن، ادامه دارد [۲]. برای بهبود انتخاب موقعیت‌یابی در طول وضعیت حمله (یعنی تیم صاحب توپ)، بازیکن‌ها باید بهترین موقعیت قابل دسترسی جهت دریافت پاس یا امتیاز گل را پیدا کنند [۱].

اولین رقابت ربات‌های فوتبال‌بالیست در سال ۱۹۹۷ در اوزاکا ژاپن برگزار شد. در این دوره از مسابقات، استراتژی تیم‌ها کاملاً عمومی و

سراسر است بود. بیشتر تیم‌ها، بازیکن‌ها را در موقعیت ثابت نگه‌داشته و بازیکنان تنها به هنگام نزدیک‌شدن توپ به آن‌ها به طرف توپ حرکت می‌کردند. از این به بعد محققان زیادی به این زمینه علاقه‌مند شدند و روش‌های پیشرفته‌تری برای کنترل این مسئله ارائه کردند. به‌طوری‌که در سال ۱۹۹۸، ربات‌های GPR2D برای اولین بار میزبان ارزیابی علم سیستم‌های چندعامله در رقابت‌های جهانی بود [۳].

در پژوهش‌های [۴-۶] برای کنترل ربات‌های فوتبال‌بالیست از یادگیری تقویتی استفاده شده‌است. پژوهش [۴] روشی است که روی تیم GPR2D تنها برای کنترل عامل صاحب توپ پیاده‌سازی شده‌است. در این روش عامل صاحب توپ برای انتخاب عمل خود براساس موقعیت هم‌تیمی‌ها و حریفان تصمیم‌گیری می‌کند. برای این منظور از الگوریتم یادگیری کیو<sup>۵</sup> استفاده شده‌است. در این پژوهش، ۷ حالت برای محیط و ۷ عمل برای عامل در نظر گرفته شده‌است. عامل‌ها در این روش طی ۱۰ بازی مقابل تیم هلیوس ۲۰۱۱ آموزش دیدند. سپس تیم در ۱۰ بازی مقابل هلیوس ۲۰۱۱ مورد آزمون قرار گرفت و به نتایج ۴ برد و ۶ مساوی رسید. قابل ذکر است که ابتدا این روش روی تیم هلیوس ۲۰۱۳ پیاده‌سازی گردید و باعث بهبود در نتایج نشد. پژوهش [۵] نیز روی تیم اکسیوم<sup>۶</sup> پیاده‌سازی شده که از یادگیری کیو برای کنترل عملکرد حمله استفاده کرده‌است. در این مقاله ۱۵ عمل برای عامل و ۱۰ حالت برای محیط در نظر گرفته شده‌است. در این پژوهش در واقع یک الگوریتم جدید برای پاس‌کاری و موقعیت‌یابی پویای عامل‌ها ارائه گردید، ولی به علت فضای حالت بزرگ با مشکل تنگنای ابعاد روبرو شد. در پژوهش [۶] از یادگیری سارسا<sup>۷</sup> تنها برای حل زیر مسئله HFO<sup>۸</sup> استفاده شده‌است. در زیر مسئله HFO به‌جای تمام بازیکنان، یک تیم مدافع و یک دروازه‌بان و یک تیم مهاجم مورد بررسی قرار می‌گیرند. در این پژوهش محیط به‌طور پیوسته در نظر گرفته شده‌است و همچنین برای عامل، ۷ عمل تعریف شده‌است.

در پژوهش [۷] به‌منظور پیش‌بینی موقعیت حریف برای هر بازیکن یک شبکه عصبی<sup>۹</sup> سه لایه با ۵ نرون<sup>۱۱</sup> در لایه ورودی و ۱۰ نرون در لایه پنهان و ۲ نرون در لایه خروجی تعریف شده‌است. وضعیت زمین در یک گام زمانی توسط یک بردار ۵ بعدی که شامل موقعیت جاری توپ، کمترین گام زمانی برای رسیدن به توپ میان همه بازیکنان و تخمین موقعیت توپ بعد از کمترین گام زمانی نمایش داده می‌شود. خروجی شبکه عصبی، موقعیت یک بازیکن حریف متناظر با وضعیت زمین را مشخص می‌کند که موقعیت حریف را حدوداً با ۱۰ میلی‌متر خطا پیش‌بینی می‌کند. ولی نتایج آزمایش‌ها با توجه به میانگین خطاها نشان می‌دهد که، روش ارائه‌شده خیلی بهتر از روش‌های پیش‌بینی سنتی عمل نمی‌کند. در [۵] نیز برای تخمین میزان امن بودن محیط از شبکه عصبی سه لایه استفاده شده‌است. در این پژوهش ۱۰ نرون در لایه خروجی و ۶ نرون در لایه پنهان و ۱ نرون برای لایه خروجی در نظر گرفته شده‌است. در نهایت این روش به عامل برای انجام پاس‌های قابل اعتمادتر و امن‌تر کمک شایانی کرد.

در بسیاری از کارهای پیشین، الگوریتم یادگیری کیو و سارسا استفاده شده‌است. در این مقاله سعی بر این است که با به‌کارگیری الگوریتم یادگیری کیو - وی<sup>۱۶</sup> عامل‌های درگیر در حمله در فوتبال ربات‌های دو بعدی آموزش ببینند. آزمایش‌هایی که در همین زمینه [۱۳] روی عملکرد الگوریتم‌های یادگیری کیو و کیو - وی انجام گرفت، نشان از برتری الگوریتم یادگیری کیو - وی نسبت به الگوریتم کیو برای حل مسئله ربات‌های فوتبالیست دارد. از آنجایی که اختلاف انواع روش‌های یادگیری تقویتی در نحوه به‌روزرسانی توابع ارزش می‌باشد، از طرفی یادگیری کیو - وی از تابع ارزش حالت برای به‌روزرسانی تابع ارزش حالت - عمل خود استفاده می‌کند، دارای نتایج بهتری نسبت به دیگر روش‌های یادگیری تقویتی است. روش‌های ارائه شده در این مقاله بر روی آخرین کد پایه در دسترس از تیم هلیوس (هلیوس ۲۰۱۳) پیاده‌سازی شده‌است. لازم به ذکر است این تیم در رقابت‌های چندساله اخیر در مسابقات جهانی روبوکاپ جزو تیم‌های اول و دوم بوده است.

ساختار این نوشتار در پنج بخش صورت گرفته است. در بخش دوم مفاهیم مربوط به یادگیری تقویتی و روش انتخاب عمل شرح داده می‌شود. همچنین توضیحاتی در ارتباط با محیط شبیه‌سازی که استفاده شده‌است، بیان شده‌است. بخش سوم شرح فرمول‌گذاری کردن مسئله آمده‌است. در بخش چهارم پیاده‌سازی و نتایج روش بیان شده‌است. و بخش آخر جهت بررسی نتایج و بیان پیشنهادات در نظر گرفته شده‌است.

## ۲- مفاهیم اولیه

در این بخش به تشریح مفاهیم و روش‌های موردنیاز برای درک بهتر مطالب این مقاله پرداخته شده‌است.

### ۲-۱- یادگیری تقویتی

یادگیری تقویتی یک شاخه از یادگیری ماشین است که عامل یادگیرنده تلاش می‌کند تا انجام یک وظیفه مشخص را در یک محیط مشخص یاد بگیرد. هدف عامل پیشینه‌کردن مجموع پاداش دریافتی در طول اجرا می‌باشد [۱۴]. در یادگیری تقویتی یک یا چند عامل با توجه به تقابلی که با محیط دارند یک عمل را در یک حالت داده‌شده انتخاب می‌کنند تا به یک سیاست بهینه برسند. سیاست بهینه یادگرفته شده برای انجام این روش باید بهترین عمل را از بین عمل‌های قابل اجرا در هر موقعیت انتخاب کند [۱۵]. ساختار کلی یادگیری تقویتی در شکل ۱ قابل مشاهده است.

یک مسئله یادگیری تقویتی توسط سه فاکتور حالت، عمل و سیگنال تقویتی فرمول‌گذاری می‌شود. یک وضعیت در محیط توسط حالت مشخص می‌شود. حالت  $S_t \in S$  وضعیتی از محیط در زمان  $t$  را تعریف می‌کند. در اینجا  $S$  مجموعه‌ای از حالت‌های ممکن است. عامل یادگیرنده در هر گام زمانی مجموعه‌ای از عمل‌هایی که می‌تواند در حالت جاری انتخاب کند، دریافت می‌کند. عامل یک عمل  $a_t \in A$  را انتخاب می‌کند که  $A$  مجموعه‌ای از عمل‌های ممکن است. انجام یک عمل عامل را به

در پژوهش [۸] از استدلال مبتنی بر مورد<sup>۱۲</sup> و طبقه‌بند بیزی<sup>۱۳</sup> جهت کنترل استراتژی حمله و دفاع استفاده شده‌است. براساس موردهای ثابتی که تعریف شده‌است، با استفاده از طبقه‌بند بیزی، موردی که بیشترین شباهت را با حالت جاری دارد انتخاب می‌کند و براساس آن عمل می‌نماید. با توجه به این که در این روش از موارد ثابت استدلال مبتنی بر مورد برای تعیین وضعیت عامل در حالت جاری استفاده می‌شود، با شرایط پویا و پیچیده فوتبال ربات‌ها سازگار نیست.

در پژوهش‌های [۷، ۹] برای کنترل ربات‌های فوتبالیست از جستجوی درختی استفاده شده‌است. در [۹] به منظور کنترل رفتار پاس‌کاری بازیکنان و موقعیت‌یابی بازیکنان درخت جستجوی پیشرو<sup>۱۴</sup> بکار رفته است. در این روش با توجه به ویژگی‌های غیر قابل پیش‌بینی از محیط، موقعیت‌یابی چندگامی و رفتار پاس‌دادن تنها در یک افق خیلی محدود امکان‌پذیر است. پژوهش [۷] روشی است که توسط تیم هلیوس ۲۰۱۳ پیاده‌سازی شده‌است. در این پژوهش برای کنترل رفتار تهاجمی عامل‌ها جستجوی اول بهترین<sup>۱۵</sup> بکار رفته است. از آنجایی که جستجوی اول بهترین استفاده شده‌است، همگرایی روش به سیاست بهینه تضمین نمی‌شود.

پژوهش [۱۰] روشی برای ارزیابی استراتژی عامل‌های فوتبال و کسب یک استراتژی مناسب با استفاده از الگوریتم ژنتیک ارائه کرده است. در این مقاله، هر عامل فوتبال دارای یک مجموعه قوانین عمل است که تنها هنگامی که توپ در ناحیه قابل ضربه‌زدن عامل است، استفاده می‌شود. این روش سرعت محاسباتی پایینی دارد و در محیط بی‌درنگ و پیچیده فوتبال ربات‌ها کارایی خوبی ندارد.

پژوهش [۱۱] نسخه بهبودیافته روش [۷] را ارائه کرده‌است که در آن از یک مدل فازی برای ارزیابی هر چه بهتر عمل‌ها استفاده شده‌است. در این پژوهش عامل استراتژی‌های خوب تیم‌های دیگر را یاد می‌گیرد و باعث بهبود توانایی پاس‌کاری عامل‌ها در مقایسه با حالت بدون این سیستم شده‌است.

روش دیگری که در [۱۰] ارائه گردید، از یک سری پارامتر برای تصمیم‌گیری جهت علامت‌گذاری و تحت‌فشار گذاشتن حریف در موقعیت حمله استفاده می‌کند. در اینجا تحت‌فشار گذاشتن، یک اقدام تاکتیکی جمعی است که در شرایط غیرتصاحبی توسط بیش از یک بازیکن انجام می‌شود. در اینجا هدف این است که فضا و زمان بازی برای تیمی که در تصاحب است بسته شود، این توسعه حرکات حمله آن‌ها را مشکل ساخته و بدست آوردن مجدد توپ را برای تیم آسان می‌سازد.

یکی از مزیت‌های یادگیری تقویتی سادگی در پیاده‌سازی است که با تعریف سه فاکتور حالت، عمل و سیگنال تقویتی برای مسئله موردنظر، به‌سادگی به حل مسئله می‌پردازد. اگرچه این برای دیگر روش‌های یادگیری ماشین ساده نیست [۱۲]. به همین علت در این مقاله برای کنترل و بهبود مهارت‌های حمله که هدف این تحقیق است از یادگیری تقویتی استفاده شده‌است. در مقاله حاضر تاکید بر روی سیستم‌های متشکل از چند عامل خودمختار است که می‌توانند در محیط‌های بی‌درنگ، نویزی، نیاز به همکاری و خصمانه با اهداف متقابل عمل کنند.

به‌روزرسانی تابع ارزش حالت - عمل به‌ترتیب در روابط (۳) و (۴) برای یادگیری کیو - وی آمده‌است:

$$V(s_t) = V(s_t) + \beta[r_t + \gamma V(s_{t+1}) - V(s_t)] \quad (3)$$

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha[r_t + \gamma V(s_{t+1}) - Q(s_t, a_t)] \quad (4)$$

در اینجا متغیرها به‌صورت زیر تعریف می‌شوند:

-  $s_t$  حالت جاری متناظر

-  $a_t$  عمل انتخاب شده در حالت  $s_t$

-  $r_t$  سیگنال تقویتی دریافتی بر اثر انتخاب عمل  $a_t$  در حالت  $s_t$

-  $s_{t+1}$  حالت بعدی

-  $\gamma$  ضریب تخفیف

-  $\alpha$  و  $\beta$  نرخ آموزش

الگوریتم ۱ شبه کد الگوریتم یادگیری کیو-وی را نمایش می‌دهد.

#### الگوریتم ۱: الگوریتم یادگیری کیو - وی

```

Initialize Q(s,a) arbitrarily
Initialize V(s) arbitrarily
Repeat (for each episode):
  Initialize S
  Repeat (for each step of episode):
    Choose  $a_{t+1}$  from  $s_{t+1}$  using policy derived from Q
    (e.g.  $\epsilon$ -greedy)
    take action a, observe  $r_{t+1}; s_{t+1}$ 
     $Q(s_t, a_t) = Q(s_t, a_t) + [r_{t+1} + \gamma V(s_{t+1}) - Q(s_t, a_t)]$ 
     $V(s_t) = V(s_t) + \beta[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$ 
     $s_t \leftarrow s_{t+1}$ 
  Until  $s_t$  is Terminal
    
```

#### ۲-۲- انتخاب عمل به روش شبه‌حریصانه

یکی از سیاست‌ها که به‌طور گسترده مورد استفاده قرار می‌گیرد روش شبه‌حریصانه<sup>۱۸</sup> می‌باشد. در این روش با احتمال  $1-\epsilon$  (  $\epsilon$  یک عدد حقیقی مثبت بین صفر و یک) عملی برگزیده می‌شود که دارای بالاترین مقدار ارزش باشد. دیگر عمل‌ها با احتمال  $\epsilon$  با شانس مساوی انتخاب می‌شوند. احتمال انتخاب هر عمل در سیاست شبه‌حریصانه در رابطه (۶) آمده‌است [۱۴].

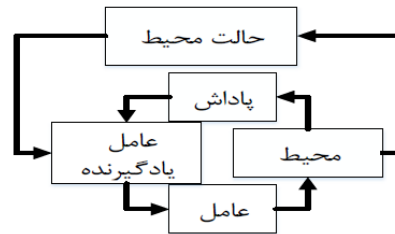
$$p(a) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{N} & a = \arg \max_{b \in A} (Q(s, b)) \\ \frac{\epsilon}{N} & otherwise \end{cases} \quad (6)$$

در رابطه بالا  $N$  تعداد عمل‌ها می‌باشد و  $\arg \max_{b \in A} (Q(s, b))$  شماره عمل با بیشینه ارزش را برمی‌گرداند.

#### ۲-۳- معرفی محیط شبیه‌ساز ربات‌های فوتبالیست

سیستم شبیه‌سازی که پیاده‌سازی‌ها و آزمون‌های این مقاله در آن صورت گرفته است، rcsoccersim15.1.1 می‌باشد. این سیستم شبیه‌ساز از سه بخش اصلی سرور فوتبال، نمایشگر فوتبال و نمایشگر مجدد بازی تشکیل شده‌است [۳]. سیستم سرور، نرم افزار اصلی در طراحی و مدل‌سازی محیط فوتبال می‌باشد. تمام اطلاعات مربوط به حسگرها، عملگرها، الگوریتم‌های ایجاد خطا و الگوهای مورد استفاده در

حالت جدید در گام زمانی بعدی منتقل می‌کند و عامل یک سیگنال تقویتی  $r_t \in R$  به عنوان نتیجه‌ای از عملش دریافت می‌کند [۱۴].



شکل ۱: تعامل عامل و محیط در الگوریتم یادگیری تقویتی [۱۶]

تقریباً همه الگوریتم‌های یادگیری تقویتی با تخمین تابع ارزش حالت (یا حالت - عمل) سروکار دارند که تابع ارزش تخمین زده‌شده توسط عامل، نحوه انتخاب عمل عامل در یک حالت داده‌شده را مشخص می‌کند. باید توجه داشت که تابع‌های ارزش با توجه به سیاست‌های خاص تعریف می‌شوند. یک سیاست،  $\pi$ ، یک نگاشت از حالت  $s \in S$  به عمل  $a \in A$  است [۱۴].

به‌طور رسمی ارزش حالت  $s$  تحت سیاست  $\pi$  به‌صورت  $V_{\pi}(s)$  نشان داده می‌شود، مقدار بازگشتی موردانتظار در صورتی که عامل با شروع از حالت  $s_t = s$  تحت سیاست  $\pi$  به دست می‌آورد را نشان می‌دهد. تابع ارزش حالت،  $V_{\pi}(s)$ ، به‌صورت رابطه (۱) تعریف می‌گردد.

$$V_{\pi}(s) = E_{\pi}[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s] \quad (1)$$

در اینجا  $E_{\pi}[\cdot]$  امید مقدار پاداش مورد انتظاری است که عامل بر اثر دنبال کردن سیاست  $\pi$  کسب می‌کند.  $\gamma \in [0,1]$  نرخ تنزل می‌باشد که اهمیت نسبی پاداش‌های بلندمدت را نسبت به پاداش‌های کوتاه‌مدت تنظیم می‌کند. متغیر  $t$  نشان‌دهنده هر گام زمانی است. همچنین باید توجه داشت که ارزش حالت پایانی در هر صورتی، صفر می‌باشد.

به‌طور مشابه ارزش عمل  $a$  در حالت  $s$  تحت سیاست  $\pi$ ، به‌صورت  $Q_{\pi}(s, a)$  نشان داده می‌شود، که مقدار بازگشتی موردانتظار با شروع از حالت  $s_t = s$  و انتخاب عمل  $a_t = a$ ، و دنبال کردن سیاست  $\pi$  می‌باشد. تابع ارزش حالت - عمل  $Q_{\pi}(s, a)$  به‌صورت رابطه (۲) تعریف می‌شود.

$$Q_{\pi}(s, a) = E_{\pi}[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a] \quad (2)$$

در این مقاله الگوریتم یادگیری کیو - وی پیاده‌سازی شده‌است و نتایج آن مورد بررسی قرار گرفته است. در ادامه این بخش نحوه عملکرد الگوریتم یادگیری کیو - وی و نحوه تخمین تابع ارزش حالت - عمل مربوط به آن آمده‌است.

#### ۲-۱-۱- یادگیری کیو - وی

یادگیری کیو - وی یک روش یادگیری تقویتی برسیاست<sup>۱۷</sup> است که تعمیمی از ایده‌های دو روش یادگیری کیو و سارسا می‌باشد. در این روش به‌جای نگهداری یک تابع ارزش، همزمان هر دو تابع ارزش حالت  $V$  و تابع ارزش حالت - عمل  $Q$  نگهداری می‌شود. که همین مسئله باعث بهبود عملکرد در روش یادگیری کیو - وی نسبت به دو روش یادگیری کیو و سارسا شده‌است [۱۴]. قاعده به‌روزرسانی تابع ارزش حالت و

برقراری ارتباط میان عامل‌ها توسط این سیستم، تولید، پردازش و بین عامل‌ها، منتقل می‌شوند.

سیستم شبیه‌ساز، دو نمایشگر بازی برخط<sup>۱۹</sup> و برون خط<sup>۲۰</sup> نیز به‌همراه دارد که به‌طور مستقیم با سرور در تقابل است و اطلاعات ایجاد شده از سوی سرور را نمایش می‌دهد. روش نمایش برخط در هنگام اجرای یک مسابقه واقعی و روش نمایش برون خط در زمان آزمون و یا مطالعه روی نحوه عملکرد یک تیم در یک بازی انجام شده مورد استفاده قرار می‌گیرد. یکی دیگر از اجزای شبیه‌ساز، نمایشگر مجدد بازی است که اطلاعات مربوط به یک بازی انجام شده را مجدداً توسط نمایشگر به روش برون خط نمایش می‌دهد.

### ۳- روند انجام پژوهش و روش پیشنهادی

با توجه به مطالعاتی که صورت گرفت و آزمایشاتی که انجام شد یادگیری تقویتی برای بهبود عملکرد حمله تیم در نظر گرفته شد. از جمله ویژگی‌های یادگیری تقویتی که کارایی مناسب آن برای کار روی محیط پیچیده و بی‌درنگ فوتبال ربات‌ها را اثبات می‌کند، می‌توان به موارد زیر اشاره نمود:

- با تعریف سه فاکتور حالت، عمل و سیگنال تقویتی قابل فرمول‌گذاری است.
- یادگیری با یک معیار عددی (سیگنال تقویتی) صورت می‌گیرد.
- به عامل توانایی کسب رفتارهای پیشرفته با دانش کم را می‌دهد.
- دارای هزینه محاسباتی پایین می‌باشد.

روند کلی کار و الگوریتمی که در هر دو روش ارائه شده به ترتیب در شکل ۲ و الگوریتم ۲ آمده‌است. با توجه به شکل ۲ و الگوریتم ۲ در ابتدا باید سه فاکتور اصلی حالت، عمل و سیگنال تقویتی و همچنین مقادیر اولیه توابع ارزش حالت و ارزش حالت - عمل مشخص گردد. سپس عامل در طول بیست بازی طی مراحل آموزش مختلف آموزش می‌بیند تا به سیاست نهایی دست پیدا کند. سیاست نهایی نحوه انتخاب عمل عامل را در حالت‌های مختلف تعیین می‌کند. نحوه تعیین فاکتورهای یادگیری تقویتی که ایده اصلی این مقاله برای بهبود عملکرد حمله در فوتبال ربات‌ها می‌باشد در زیر بخش‌های بعدی آمده‌است.

#### الگوریتم ۲: الگوریتم پیشنهادی

- ۱- گسسته‌سازی حالت‌ها و عمل‌ها
- ۲- تعیین سیاست دلخواه برای هر حالت و مقداردهی پارامترها
- ۳- شروع حلقه ۱ (به ازای هر بازی)
- ۴- شروع حلقه ۲ (به ازای هر مرحله آموزش)
- ۵- انتخاب عمل برای هر حالت
- ۶- به‌کارگیری عمل‌ها در محیط
- ۷- دریافت ادراک از محیط
- ۸- تعیین حالت و سیگنال تقویتی

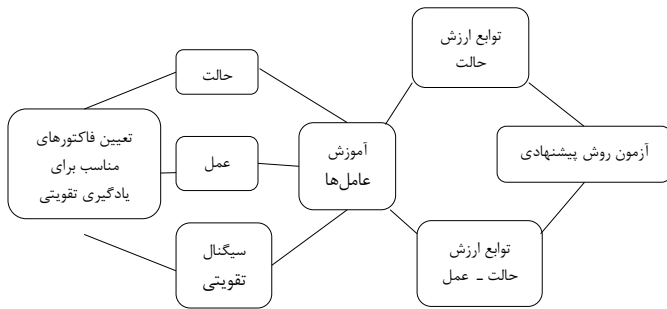
۹- به‌روزرسانی توابع ارزش حالت و ارزش حالت - عمل

۱۰- اگر در حالت پایانی قرار ندارد یا توپ از منطقه جریمه خارج نشده

برو به گام ۴

۱۱- اگر بازی به پایان نرسیده برو به گام ۳

۱۲- برگرداندن سیاست



شکل ۲: روند انجام پژوهش

#### ۳-۱- فرمول‌گذاری کردن مسئله

جهت پیاده‌سازی الگوریتم یادگیری کیو - وی در ابتدا لازم است تا فضای حالتی که عامل قادر به درک آن است گسسته‌سازی شود. از آنجایی که حالت‌هایی که عامل صاحب توپ و عامل‌های بدون توپ درک می‌کنند با هم متفاوت است، عامل‌ها در دو دسته عامل صاحب توپ و عامل بدون توپ در نظر گرفته شد. جهت گسسته‌سازی حالت‌ها با توجه به اطلاعاتی که عامل از محیط اطراف دارد، برای عامل صاحب توپ یک بردار ۶ بیتی و برای عامل بدون توپ یک بردار ۳ بیتی جهت تعیین ویژگی‌های محیط برای حالت متناظر عامل در نظر گرفته شده‌است. در زیر شرح کاملی از ویژگی‌های در نظر گرفته شده و نحوه گسسته‌سازی فضای حالت آمده‌است.

بردار ۶ بیتی که جهت گسسته‌سازی حالت‌های عامل صاحب توپ در نظر گرفته شد به شرح زیر می‌باشد:

- بیت اول: نشان‌دهنده قابلیت پاس‌دادن توسط عامل صاحب توپ در حالت جاری است.

- بیت دوم: نشان‌دهنده حضور یا عدم‌حضور عامل صاحب توپ روبروی دروازه است.

- بیت سوم: نشان‌دهنده وضعیت نزدیک‌ترین بازیکن حریف است، یعنی بیان می‌کند موقعیت حریف پشت یا جلو بازیکن است.

- بیت چهارم: اگر فاصله نزدیک‌ترین بازیکن حریف تا بازیکن صاحب توپ کمتر از ۲ متر باشد این بیت ۱ می‌شود.

- بیت پنجم: اگر فاصله نزدیک‌ترین بازیکن حریف تا بازیکن صاحب توپ بین ۲ تا ۴ متر باشد این بیت ۱ می‌شود.

- بیت ششم: اگر فاصله نزدیک‌ترین بازیکن حریف تا بازیکن صاحب توپ بیشتر از ۴ متر باشد این بیت ۱ می‌شود.

مورد آزمون قرار گرفت. شرح مفصلی از روش‌های پیاده‌سازی شده در ادامه آمده‌است.



شکل ۳: نمایش منطقه جریمه، ناحیه آموزش عامل صاحب توپ

### ۳-۱-۱- توزیع یکسان سیگنال تقویتی

در اولین روشی که پیاده‌سازی شد تمام عامل‌های یادگیری تقویتی به‌طور یکسان سیگنال تقویتی دریافت می‌کنند. در واقع به‌هنگام رسیدن سیگنال پاداش یا جریمه تمام عامل‌های درگیر در حمله سیگنال تقویتی را به‌طور یکسان دریافت می‌کنند. با ورود توپ به منطقه جریمه آموزش عامل شروع می‌شود و یک مرحله آموزش آغاز می‌شود و اگر هر یک از سه حالت پایانی زیر اتفاق بیفتد مرحله به پایان می‌رسد و عامل‌ها سیگنال تقویتی به شرح زیر را دریافت می‌کنند.

- اگر بر اثر عمل توسط عامل، توپ از منطقه جریمه خارج شود، عامل جریمه ۵- دریافت می‌کند.
- اگر بر اثر عمل عامل، توپ از دست تیم خودی خارج شود و به دست تیم حریف برسد، عامل جریمه ۵۰- را دریافت خواهد کرد.
- اگر عامل در حالتی قرار بگیرد که روبروی دروازه باشد و حریف در مقابل آن قرار نداشته باشد و قابلیت پاس دادن هم وجود نداشته باشد (حالت‌های ۶، ۷ و ۸). در این صورت عامل عمل kick را انتخاب می‌کند و پاداش ۱۰۰ را دریافت خواهد کرد. در این مورد جهت همکاری بهتر عامل‌ها و پاس‌کاری بیشتر بین عامل‌ها قابلیت پاس‌دادن مورد بررسی قرار گرفته‌است.

به‌علاوه از آنجایی‌که در یادگیری تقویتی عامل به‌دنبال بیشینه‌کردن پاداش دریافتی در طول اجرا می‌باشد، برای اینکه عامل در گام زمانی کمتری به گل برسد در هر گام زمانی جریمه ۵/۰- برای عامل در نظر گرفته شده‌است.

عامل‌ها در طول ۲۰ بازی مقابل تیم هلیوس ۲۰۱۳ آموزش دیدند تا تابع Q عامل‌ها به سیاست نهایی همگرا شود. جهت همگرایی سریع‌تر

از آنجایی‌که تنها یکی از ۳ بیت آخر به‌طور همزمان یک می‌شود، ۲۴ حالت به دست می‌آید، که بر حسب مقادیر از ۰ تا ۲۳ شماره‌گذاری می‌شوند.

بردار ۳ بیتی که جهت گسسته‌سازی حالت‌های عامل‌های بدون توپ در نظر گرفته شد به شرح زیر می‌باشد:

- بیت اول: بیانگر این است که بازیکن سریع‌ترین بازیکن به توپ (یعنی نزدیک‌ترین بازیکن به بازیکن صاحب توپ) هست یا خیر
  - بیت دوم: اگر بازیکن توانایی دریافت پاس از بازیکن هم‌تیمی خود را داشته باشد این بیت یک می‌شود.
  - بیت سوم: نشان دهنده وضعیت نزدیک‌ترین بازیکن حریف است که پشت یا جلو بازیکن قرار دارد.
- در نتیجه ۸ حالت برای عامل بدون توپ به وجود می‌آید که از ۰ تا ۷ شماره‌گذاری می‌شود.

مجموعه عمل‌های عامل صاحب توپ به شرح زیر است:

- Slow\_Forward: حرکت آهسته همراه با توپ
  - Fast\_Forward: حرکت سریع همراه با توپ
  - Pass: پاس‌دادن توپ به نزدیک‌ترین بازیکن هم‌تیمی
  - Dribble: تلاش برای دریبل کردن بازیکن حریف
  - Hold\_ball: نگه‌داشتن توپ
  - Kick: ضربه مستقیم به توپ به سمت دروازه حریف
- در حالت‌هایی که بیت سوم صفر باشد یعنی نزدیک‌ترین بازیکن حریف پشت سر بازیکن صاحب توپ قرار دارد عمل Dribble از عمل‌های عامل حذف شده‌است. زیرا عمل دریبل کردن حریف تنها هنگامی‌که بازیکن حریف جلوی توپ قرار دارد معنا پیدا می‌کند. بنابراین برای عامل صاحب توپ ۶ عمل درحالتی که نزدیک‌ترین بازیکن حریف جلوی توپ قرار دارد و ۵ عمل درحالتی که نزدیک‌ترین بازیکن حریف پشت سر توپ قرار دارد در نظر گرفته شده‌است.

مجموعه عمل‌های عامل‌های بدون توپ به شرح زیر است:

- Get\_Ball: موقعیت‌یابی در بهترین موقعیت جهت به‌دست‌آوردن توپ
- Get\_Open: موقعیت‌یابی پایه تیم

از آنجایی‌که قسمت اصلی حمله در منطقه جریمه صورت می‌گیرد (مستطیل قرمز رنگی که در شکل ۳ مشاهده می‌کنید)، در این شبیه‌سازی نیز عامل تنها با ورود توپ به این محوطه آموزش می‌بیند و از روش یادگیری تقویتی ارائه‌شده استفاده می‌کند.

در این مقاله تمام پیاده‌سازی‌ها و بهبودها بر روی کد پایه هلیوس ۲۰۱۳ انجام شده‌است. روش یادگیری کیو - وی با انتخاب عمل به روش شبه‌حریصانه پیاده‌سازی شده‌است. با توجه به مطالعاتی که روی پژوهش‌های پیشین انجام شد، برای پارامترهای مربوط به الگوریتم یادگیری کیو - وی  $\gamma$ ،  $\alpha$  و  $\beta$  به ترتیب مقادیر ثابت ۰/۸، ۰/۵ و ۰/۵ و برای پارامتر  $\epsilon$  در روش انتخاب عمل شبه‌حریصانه مقدار ۰/۲ در نظر گرفته شد. بازیکنان حمله با دو رویکرد مختلف، به‌طور جداگانه آموزش دیدند و عملکرد آن‌ها نسبت به یکدیگر و تیم پایه هلیوس ۲۰۱۳

آزمون روش‌ها به صورت جدول ۱ به دست آمد. جهت مقایسه بهتر بین نتایج، ستون امتیاز آمده است. برای محاسبه امتیاز هر تیم در طول ۶۰ بازی، امتیاز ۳ برای برد، امتیاز ۱ برای مساوی و امتیاز ۰ برای باخت در نظر گرفته شده است و مجموع آن‌ها محاسبه گردیده است.

جدول ۱: نتایج بازی روش‌های ارائه شده مقابل تیم هلیوس ۲۰۱۳

تیم‌ها	برد	مساوی	امتیاز
توزیع یکسان سیگنال تقویتی هلیوس	۳۴۹	۸	۱۷۱۵۵
تقسیم سیگنال تقویتی براساس خبرگی هلیوس	۳۵۱	۶	۱۵۱۵۹

نتایج پیاده‌سازی نشان می‌دهد که روش‌های ارائه شده نتایج بهتری در مقابل تیمی که آموزش ندیده است، دارند و باعث بهبود عملکرد کلی تیم شده است. همچنین روش تقسیم سیگنال تقویتی براساس خبرگی باعث می‌شود، عامل سریع‌تر به هدف و سیاست نهایی برسد و باعث بهبود در سرعت آموزش عامل‌ها شده است. به علاوه روش ارائه شده باعث افزایش دانش به دست آمده در طول یادگیری به خصوص در اوایل یادگیری و در نتیجه افزایش سرعت یادگیری می‌شود.

به منظور اینکه صحت نتایج بالا یعنی سرعت بیشتر آموزش در روش تقسیم سیگنال تقویتی براساس خبرگی نسبت به روش توزیع یکسان سیگنال تقویتی نشان داده شود، آموزش بجای ۲۰ بازی با ۱۵ بازی تکرار شد. پس از ۱۵ آموزش در مقابل تیم هلیوس ۲۰۱۳، عملکرد هر دو تیم مقابل هم با ۶۰ بازی مورد آزمون قرار گرفت. نتایج این ۶۰ آزمون در جدول ۲ آمده است. جدول ۲ نشان‌دهنده نتایج آزمون ۱ است، یعنی حالتی که در آن تیم‌ها طی ۲۰ بازی آموزش دیدند. جدول ۲ همچنین نشان‌دهنده نتایج آزمون ۲ است، یعنی حالتی که تیم‌ها طی ۱۵ بازی آموزش دیدند. در هر دو آزمون، دو تیم آموزش دیده در حالت توزیع یکسان سیگنال تقویتی و حالت تقسیم سیگنال تقویتی براساس خبرگی مقابل هم در ۶۰ بازی مورد آزمون قرار گرفتند.

نتایج جدول ۲ نشان‌دهنده همگرایی سریع‌تر روش تقسیم سیگنال تقویتی براساس خبرگی به سیاست نهایی نسبت به روش توزیع یکسان سیگنال تقویتی می‌باشد. همان‌طور که پیداست روش تقسیم سیگنال تقویتی براساس خبرگی موفق به کسب امتیاز بیشتری نسبت به روش توزیع یکسان سیگنال تقویتی در ۱۵ آموزش شد که نشانگر عملکرد بهتر و همگرایی سریع‌تر به سیاست نهایی در روش تقسیم سیگنال تقویتی براساس خبرگی می‌باشد.

جدول ۲: نتایج ۶۰ تست روش تقسیم سیگنال تقویتی براساس خبرگی در مقابل روش توزیع یکسان سیگنال تقویتی

تیم‌ها	برد	مساوی	باخت	امتیاز
آزمون ۱ (آموزش تیم‌ها طی ۲۰ بازی)	۲۲	۲۱	۱۷	۸۷
آزمون ۲ (آموزش تیم‌ها طی ۱۵ بازی)	۴۴	۱۰	۶	۱۴۲

برای مقایسه بهتر روش‌ها نمودار شکل ۴ آمده است که مجموع امتیازات کسب شده توسط تیم آموزش دیده با هر دو روش توزیع یکسان سیگنال تقویتی و توزیع سیگنال تقویتی براساس خبرگی، که در این

آموزش به سیاست نهایی در انتهای هر بازی از تابع ارزش حالت - عمل به دست آمده توسط هر یک از بازیکنان حمله در حالت بدون توپ میانگین گرفته و به عنوان تابع ارزش حالت - عمل اولیه عامل در بازی بعد در نظر گرفته شد.

### ۳-۱-۲- تقسیم سیگنال تقویتی براساس خبرگی

دانش در عامل یادگیری تقویتی نمایانگر این است که عامل کاوش بیشتری روی محیط داشته است و تجربیات بیشتری کسب کرده است. در واقع دانش بیشتر نشان‌دهنده این است که عامل جفت حالت - عمل بیشتری را ملاقات کرده است. معیاری که در این مقاله جهت بررسی میزان خبرگی عامل‌ها در نظر گرفته شده است، به صورت رابطه (۷) است که در آن دانش عامل با محاسبه اختلاف میان بیشینه و کمینه مقادیر در تابع ارزش حالت - عمل عامل در یک حالت داده شده، سنجیده می‌شود. با توجه به این رابطه در واقع عاملی که اختلاف میان بیشینه و کمینه مقادیر در تابع ارزش حالت - عملش کمتر باشد دانش کمتری دارد.

$$C_i^K = \max(Q_i(s_t, a)) - \min(Q_i(s_t, a)) \quad (7)$$

$$i = 1, \dots, n$$

ایده اصلی در این روش بر این اساس بود که، به هنگامی که عملکرد یک سیستم چندعامله به درستی نمی‌باشد تمام عامل‌ها مقصر نمی‌باشند، و عامل با دانش کمتر احتمال خطای بیشتری دارد. بر همین اساس جهت تقسیم بهتر سیگنال تقویتی بین عامل‌ها و آموزش بهتر بازیکنان به هنگام رسیدن جریمه ۵۰- در طول بازی‌ها هر عامل براساس دانشی که دارد مطابق با رابطه (۸) جریمه دریافت می‌کند و بقیه سیگنال‌های تقویتی مشابه با روش قبل توزیع می‌گردد. به واقع در این روش عاملی که دانش کمتری دارد جریمه بیشتری نسبت به عاملی که دانش بیشتری دارد دریافت می‌کند. این امر باعث می‌شود عامل‌ها تلاش بیشتری برای کسب دانش پیدا کنند و با سرعت بیشتری به سیاست نهایی همگرا شوند. به خصوص در ابتدای آموزش که عامل‌ها دانش کمتری دارند باعث افزایش سرعت آموزش می‌گردد. عامل‌ها در این روش نیز در طول ۲۰ بازی در مقابل تیم هلیوس ۲۰۱۳ آموزش دیدند تا تابع Q به سیاست نهایی همگرا گردد.

$$P_i = \left( 1 - \frac{\max(Q_i(s_t, a)) - \min(Q_i(s_t, a))}{\max(\max(Q_j(s_t, a)) - \min(Q_j(s_t, a)))} \right) \quad (8)$$

$$* (-50)$$

$$i = 1, \dots, n \quad j = 1, \dots, n$$

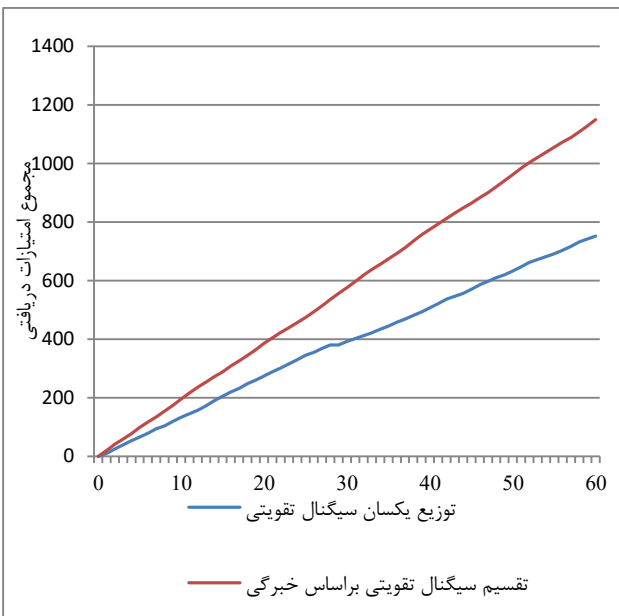
### ۴- نتایج پیاده‌سازی

برای آزمون روش‌هایی که در بالا شرح داده شد ۶۰ بازی به طور جداگانه بین روش توزیع یکسان سیگنال تقویتی بین عامل‌ها و تیم هلیوس ۲۰۱۳ و ۶۰ بازی میان روش تقسیم سیگنال تقویتی بین عامل‌ها براساس خبرگی و تیم هلیوس ۲۰۱۳ برگزار شد. نتایج بازی‌های انجام شده جهت

نتایج آزمون‌ها حاکی از آن است که روش‌های ارائه‌شده توانایی مقابله با تیم‌های برتر جهانی و حریفان مختلف را دارد و بهبود خوبی روی عملکرد تیم داشته‌است. همچنین عامل‌ها توانایی انتخاب بهترین عمل در هر حالت را با کمک سیاستی که در طول آموزش کسب کردند، دارند.

جدول ۳: نتایج بازی روش‌های ارائه‌شده مقابل تیم ربات‌ایگل ۲۰۱۳

تیم‌ها	برد	مساوی	امتیاز
توزیع یکسان سیگنال تقویتی/رایت‌ایگل	۰/۵۹	۱	۱/۱۷۸
تقسیم سیگنال تقویتی براساس خبرگی/رایت‌ایگل	۰/۶۰	۰	۰/۱۸۰



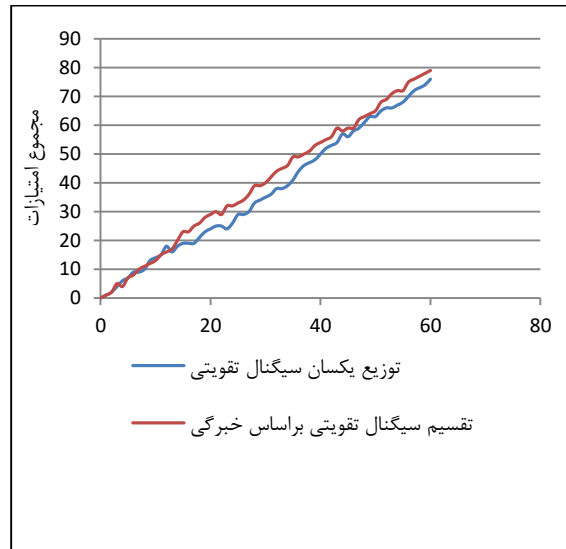
شکل ۶: مقایسه مجموع امتیازات کسب‌شده روش‌های ارائه‌شده در

۶۰ بازی مقابل تیم ربات‌ایگل ۲۰۱۳

برای مقایسه بهتر نتایج از لحاظ آماری، از آزمون فریدمن جهت رتبه‌بندی روش‌ها استفاده شده‌است. آزمون فریدمن که به آزمون تحلیل واریانس معروف است، یکی از آزمون‌های آماری است که از آن برای مقایسه میانگین رتبه‌های K متغیر یا گروه استفاده می‌شود. در اینجا با اعمال آزمون فریدمن بر روی نتایج آزمون‌ها به رتبه‌بندی روش‌ها پرداخته شده‌است. برای این منظور تعداد گل‌هایی که هر تیم در طول آزمون زده‌اند، ملاک رتبه‌بندی قرار گرفته‌است. نتایج آزمون که با استفاده از نرم‌افزار spss محاسبه شده‌است، به شرح شکل ۷ است.

در شکل ۷ جدول Rank میانگین رتبه‌های هر تیم را نشان می‌دهد. جدول Test Statistics محتوی نتایج اصلی آزمون است که مقدار آماره کی‌دو با یک درجه آزادی را نشان می‌دهد. همچنین سطح معنی‌داری آزمون را برابر با مقدار صفر نشان می‌دهد که بیانگر این است که بین میانگین رتبه‌ها اختلاف معناداری وجود دارد. با توجه به جداول Rank در شکل ۷ هر دو روش ارائه‌شده امتیاز بیشتری مقابل حریفان خود

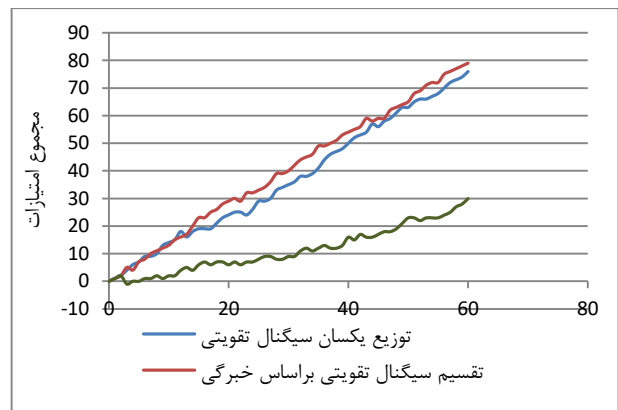
مقاله ارائه گردید، در مقابل تیم هلیوس ۲۰۱۳ (بعد از ۲۰ بازی آموزشی) را نشان می‌دهد. در اینجا منظور از مجموع امتیارات، اختلاف گل زده توسط روش ارائه‌شده از گل خورده از تیم هلیوس ۲۰۱۳ در طول ۶۰ بازی می‌باشد.



شکل ۴: مقایسه مجموع امتیازات کسب‌شده روش‌های ارائه‌شده در

۶۰ بازی مقابل تیم هلیوس ۲۰۱۳

چندین مرحله شبیه‌سازی انجام شد. در اولین مرحله شبیه‌سازی‌ها ابتدا عامل‌های صاحب توپ آموزش دیدند [۱۳]. نتایج مقایسه با بهترین نتایج از این مرحله از شبیه‌سازی در نمودار شکل ۵ قابل مشاهده است. همچنین جهت انجام آزمون‌های بیشتر روی روش‌های ارائه‌شده دو روش پیشنهادی ۶۰ بار مقابل کد پایه آزادشده توسط تیم ربات‌ایگل<sup>۲۱</sup> در سال ۲۰۱۳ قرار گرفت. تیم ربات‌ایگل در روبرو کاپ‌های سال‌های اخیر جز تیم‌های برتر جهان بوده‌است. نتایج حاصل از این ۶۰ بازی به ترتیب در جدول ۳ و نمودار شکل ۶ قابل مشاهده است.



شکل ۵: مقایسه مجموع امتیازات کسب‌شده در ۶۰ بازی مقابل تیم

هلیوس ۲۰۱۳



همچنین باعث بهبود عملکرد تیم در مقابله با حریفان مختلف شد. به‌علاوه در روش دوم که عامل‌ها براساس سطح خبرگی خود سیگنال تقویتی دریافت کردند، همگرایی به سیاست نهایی سریع‌تر صورت گرفت و باعث افزایش سرعت در آموزش عامل‌ها گردید. با توجه به پیاده‌سازی‌ها روش دوم یعنی تقسیم سیگنال تقویتی براساس خبرگی، عملکرد بهتری نسبت به روش اول یعنی توزیع یکسان سیگنال تقویتی داشت. به‌طور کلی روش تقسیم سیگنال تقویتی براساس خبرگی باعث افزایش سرعت آموزش عامل‌ها در طول مراحل یادگیری و بهبود عملکرد تیم نسبت به تیمی که عامل‌های حمله آن آموزش ندیده‌اند، شد.

روشی که در این مقاله ارائه شد تنها برای بهبود عملکرد بازیکنان حمله در حساس‌ترین منطقه حمله یعنی منطقه نزدیک به دروازه ارائه شد. برای کارهای آینده می‌توان مسئله یادگیری تقویتی را به‌نحوی برای آموزش کل بازیکن‌ها و در کل زمین بازی فوتبال تعریف کرد. همچنین می‌توان با تعریف جزئی‌تر مسئله روی رفتارهای سطح پایین عامل‌ها، روش ارائه‌شده را برای کنترل رفتار بازیکنان در فوتبال شبیه‌ساز سه‌بعدی تعمیم داد.

## مراجع

- [1] F. Almeida, N. Lau, and L. P. Reis, "A Survey on Coordination Methodologies for Simulated Robotic Soccer Teams," RoboCup Symposium, 2010.
- [2] M. Alavi, M. F. Tarazkouhi, A. Azaran, A. Nouri, S. Zolfaghari, and H. R. S. Boroujeni, "Robocup 2012- Soccer Simulation League 2D Soccer Simulation Riton," Robot Soccer World Cup, Springer Berlin Heidelberg, 2013.
- [3] M. Chen et al., (2003) RoboCup Soccer Server for Soccer Server Tersion cefc and later, [Online], Available: <http://wwfc.cs.virginia.edu/documentation/manual.pdf> [jol. 11, 2015]
- [4] J. R. F. Neri, M. R. Zатели, C. H. F. dos Santos, and J. A. Fabro, "A Proposal of QLearning to Control the Attack of a 2D Robot Soccer Simulation Team," Robotics Symposium and Latin American Robotics Symposium (SBR-LARS), pp. 174-178, 2012.
- [5] M. Ghazanfari, S. O. Shirshorshidi, and F. Samsampour, "Axiom 2013 Team Description Paper," Robot Soccer World Cup, Springer Berlin Heidelberg, vol. 8371, 2014.
- [6] S. Kalyanakrishnan, Y. Liu, and P. Stone, "Half field offense in RoboCup soccer: A multiagent reinforcement learning case study," Robot Soccer World Cup, Springer Berlin Heidelberg, vol. 4434, pp. 72-85, 2008.
- [7] H. Akiyama, T. Nakashima, and K. Yamashita, "Helios2013 team description paper," Robot Soccer World Cup, Springer Berlin Heidelberg, vol. 8371, 2014.
- [8] T. Sirinivasan, K. Aarathi, S. A. Meenakshi, and M. Kausalya, "Cbrrobosoc: An efficient planning strategy for robotic soccer using case based reasoning," International Conference on Computational Intelligence for Modeling Control and Automation, and International Conference on Intelligent Agents, Web Technologies and Internet Commerce, pp. 113-119, 2006.
- [9] A. Bai, H. Zhang, G. Lu, M. Jiang, and X. Chen, "WrightEagle 2D Soccer Simulation Team Description," Robot Soccer World Cup, Springer Berlin Heidelberg, vol. 7500, 2013.
- [10] S. Marian, D. Luca, B. Sarac, and O. Cotarlea, "OXSY 2014 Team Description," Robot Soccer World Cup, Springer Berlin Heidelberg, 2015.
- [11] H. Akiyama, T. Nakashima, and K. Yamashita, "HELIOS2014 Team Description Paper," Robot Soccer World Cup, Springer Berlin Heidelberg, 2015.

کسب نمودند. براین اساس نتایج آزمون فریدمن بیانگر عملکرد بهتر روش‌های ارائه‌شده می‌باشد.

### Friedman Test

Ranks	
	Mean Rank
تقسیم سیگنال تقویتی براساس خبرگی	1.88
هلیوس	1.13

### Test Statistics<sup>a</sup>

N	60
Chi-Square	38.208
df	1
Asymp. Sig.	.000

(ب)

### Friedman Test

Ranks	
	Mean Rank
تقسیم سیگنال تقویتی براساس خبرگی	2.00
رایت‌ایگل	1.00

### Test Statistics<sup>a</sup>

N	60
Chi-Square	60.000
df	1
Asymp. Sig.	.000

(د)

### Friedman Test

Ranks	
	Mean Rank
توزیع یکسان سیگنال تقویتی	1.88
هلیوس	1.12

### Test Statistics<sup>a</sup>

N	60
Chi-Square	40.692
df	1
Asymp. Sig.	.000

(الف)

### Friedman Test

Ranks	
	Mean Rank
توزیع یکسان سیگنال تقویتی	1.99
رایت‌ایگل	1.01

### Test Statistics<sup>a</sup>

N	60
Chi-Square	59.000
df	1
Asymp. Sig.	.000

(ج)

شکل ۷: نتایج آزمون فریدمن، الف) روش توزیع یکسان سیگنال تقویتی مقابل تیم هلیوس، ب) روش تقسیم سیگنال تقویتی براساس خبرگی مقابل تیم هلیوس، ج) روش توزیع یکسان سیگنال تقویتی مقابل تیم رایت‌ایگل، د) روش تقسیم سیگنال تقویتی براساس خبرگی مقابل تیم رایت‌ایگل

## ۵- نتیجه‌گیری

در این مقاله روشی برای بهبود عملکرد حمله در تیم شبیه‌ساز فوتبال ربات‌ها با استفاده از یادگیری تقویتی ارائه‌شد. به‌همین منظور الگوریتم یادگیری کیو - وی روی تیم ربات‌های فوتبالیست پیاده‌سازی شد و به آموزش عامل‌های حمله پرداخته شد. برای آموزش عامل‌ها از دو روش مختلف جهت تقسیم سیگنال تقویتی بین عامل‌ها استفاده شد. در روش اول سیگنال تقویتی به‌طور یکسان بین عامل‌ها توزیع گردید و در روش دوم از ایده تقسیم سیگنال تقویتی برحسب میزان خبرگی عامل‌ها بهره گرفته‌شد. براساس این ایده عامل با دانش کمتر جریمه بیشتری نسبت به عامل با دانش بیشتر دریافت می‌کند. نتایج پیاده‌سازی‌ها نشان داد که هر دو روشی که برای آموزش عامل‌های حمله در این مقاله پیاده‌سازی شد نتایج بهتری نسبت به حالتی که عامل‌ها آموزش ندیده‌اند، داشت.

برای یک سیستم قدرت به هم پیوسته شامل SMES» مجله مهندسی برق دانشگاه تبریز، جلد ۴۷، شماره ۲، صفحات ۳۸۱-۳۹۰، تابستان ۱۳۹۶.

[۱۶] مریم رضانیان لنگرودی، سیدمازیار میرحسینی مقدم، بهنام علیزاده، «استفاده از روش یادگیری رقابتی برای قیمت‌دهی استراتژیک شرکت‌های تولید براساس LMP در بازار برق»، مجله مهندسی برق دانشگاه تبریز، جلد ۴۷، شماره ۲، صفحات ۵۳۷-۵۴۹، تابستان ۱۳۹۶.

[12] M. Yoon, *Developing basic soccer skills using reinforcement learning for the RoboCup Small Size League*, Master Thesis, Stellenbosch University, pp.11, March 2015.

[۱۳] مینا خاکسار، ولی درهمی و مهدی رضائیان، «بهبود عملکرد حمله در تیم ربات‌های شبیه‌ساز فوتبال با استفاده از یادگیری تقویتی»، دومین کنفرانس محاسبات تکاملی و هوش جمعی، دانشگاه شهید باهنر، کرمان، اسفند ۹۵.

[14] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*, MIT press Cambridge, 1998.

[۱۵] عادل اکبری مجد، حسین شایقی، حمید محمدنژاد، عبدالله یونسی، «کنترل کننده مقاوم تطبیقی بار فرانسن مبتنی بر یادگیری تقویتی

## زیرنویس‌ها

<sup>۱۲</sup> Case Based Reasoning

<sup>۱۳</sup> Bayesian Classifier

<sup>۱۴</sup> Forward Search Tree

<sup>۱۵</sup> Best First Search

<sup>۱۶</sup> QV-Learning

<sup>۱۷</sup> On-Policy

<sup>۱۸</sup>  $\epsilon$ -Greedy

<sup>۱۹</sup> Online

<sup>۲۰</sup> Offline

<sup>۲۱</sup> WrightEagle

<sup>۱</sup> RoboCup

<sup>۲</sup> Positioning

<sup>۳</sup> Offensive

<sup>۴</sup> Defensive

<sup>۵</sup> Marking

<sup>۶</sup> Q - Learning

<sup>۷</sup> Axiom

<sup>۸</sup> Sarsa

<sup>۹</sup> Half Field Offensive

<sup>۱۰</sup> Neural Network

<sup>۱۱</sup> Neuron