



بسم الله الرحمن الرحيم



دانشگاه صنعتی اصفهان

دانشکده مهندسی برق و کامپیوتر

**بهبود کیفیت و سرعت یادگیری در سیستم‌های چندعامله با استفاده از  
معیار جدید خبرگی و انتگرال فازی**

پایان‌نامه کارشناسی ارشد مهندسی کامپیوتر – هوش مصنوعی و رباتیک

داریوش حسن‌پورآده

استاد راهنما

دکتر مازیار پالهنک

پاییز ۱۳۹۵



دانشگاه صنعتی اصفهان

دانشکده مهندسی برق و کامپیوتر

پایان نامه کارشناسی ارشد رشته مهندسی کامپیوتر – هوش مصنوعی و رباتیک آقای

داریوش حسن پور آده

تحت عنوان

بهبود کیفیت و سرعت یادگیری در سیستم‌های چندعامله با استفاده از

معیار جدید خبرگی و انتگرال فازی

در تاریخ ۱۳۹۵/۱۰/۲۰ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهایی قرار گرفت:

دکتر مازیار پالهنک

۱- استاد راهنمای پایان نامه

دکتر عبدالرضا میرزایی

۳- استاد داور

دکتر محمد حسین منشئی

۴- استاد داور

دکتر محمد رضا تابان

سرپرست تحصیلات تکمیلی دانشکده

**تشکر و قدردانی**

پروردگار منّان را سپاسگزارم .....

کلیه حقوق مادی مترتب بر نتایج مطالعات،  
ابتکارات و نوآوری‌های ناشی از تحقیق  
موضوع این پایان‌نامه متعلق به دانشگاه  
صنعتی اصفهان است.

دلتنگی های آدمی را باد ترانه ای می خواند  
رویا هایش را آسمان پر ستاره نادیده می گیرد  
و هر دانه ی برفی به اشکی نریخته می ماند.  
سکوت سرشار از سخنان ناگفته است؛  
از حرکات ناکرده،  
اعتراف به عشق های نهان،  
و شگفتی های به زبان نیامده،  
در این سکوت حقیقت ما نهفته است؛  
حقیقت تو و من.

برای تو و خویش  
چشمانی آرزو می کنم،  
که چراغ ها و نشانه ها را در ظلمات مان ببیند.  
گوشی،  
که صداها و شناسه ها را در بیهوشی مان بشنود.  
برای تو و خویش،  
روحي،  
که این همه را در خود گیرد و بپذیرد.  
و زبانی  
که در صداقت خود ما را از خاموشی خویش بیرون کشد،  
و بگذارد از آن چیزها که در بندها کشیده است، سخن بگوییم.

پنجه درافکنده ایم با دست هایمان  
به جای رها شدن  
سنگین سنگین بر دوش می کشیم  
بار دیگران را  
به جای همراهی کردن شان!  
عشق ما نیازمند رهایی است نه تصاحب  
در راه خویش ایثار باید نه انجام وظیفه...

بی اعتمادی دری است  
خودستایی، چفت و بست غرور است  
و تهی دستی، دیوار است و لولا است  
زندانی را که در آن محبوس رای خویش ایم  
دلتنگی مان را برای آزادی و دلخواه دیگران بودن  
از رخنه هایش تنفس می کنیم...

# فهرست مطالب

عنوان	صفحه
فهرست مطالب	هشت
فهرست تصاویر	یازده
فهرست جداول	سیزده
چکیده	۱
<b>فصل اول : مقدمه</b>	۲
۱-۱ یادگیری مشارکتی در سیستم های چند عامله	۳
۲-۱ اهداف و نوآوری های پایان نامه	۵
۳-۱ ساختار پایان نامه	۶
<b>فصل دوم : مرور کارهای پیشین</b>	۷
۱-۲ مقدمه	۷
۲-۲ اشتراک گذاری اطلاعات	۸
۳-۲ یادگیری مشترک	۸
۴-۲ تقلید	۹
۵-۲ حافظه جمعی	۹
۶-۲ پند	۱۰
۷-۲ یادگیری مشارکتی بر مبنای خبرگی	۱۰
۸-۲ یادگیری مشارکتی بر مبنای تخته سیاه	۱۲
۹-۲ یادگیری تقویتی تعاملی	۱۲
۱۰-۲ یادگیری مشارکتی بر مبنای پختگی سیاست	۱۴
۱۱-۲ یادگیری مشارکتی بر مبنای خبرگی چند معیاری	۱۴
۱۲-۲ تسریع یادگیری مشارکتی با بهره گیری از کوتاه ترین فاصله تجربه شده	۱۵
۱۳-۲ نتیجه گیری	۱۶
<b>فصل سوم : مفاهیم علمی پیش نیاز پایان نامه</b>	۱۷
۱-۳ مقدمه	۱۷
۲-۳ یادگیری تقویتی	۱۸



۱۸	۳-۳ روش‌های انتخاب عمل
۱۹	۳-۳-۱ ε-حریصانه
۱۹	۳-۳-۲ بولتزمن
۲۰	۴-۳ الگوریتم مورد مقایسه با روش پیشنهادی
۲۰	۳-۴-۱ معیار کوتاه‌ترین فاصله تجربه‌شده
۲۱	۳-۴-۲ شوک
۲۱	۵-۳ محیط‌های آزمایش
۲۲	۳-۵-۱ محیط پلکان مارپیچ
۲۲	۳-۵-۲ محیط صید و صیاد
۲۵	۶-۳ معیارهای ارزیابی
۲۵	۷-۳ اندازه‌گیری و انتگرال فازی
۲۸	۸-۳ نتیجه‌گیری

#### فصل چهارم: روش پیشنهادی

۲۹	۱-۴ مقدمه
۳۰	۲-۴ معیار خبرگی - ماتریس ارجاع و خاطره
۳۳	۳-۴ یادگیری مشارکتی $Q$ با استفاده از ماتریس ارجاع و انتگرال فازی
۳۴	۴-۳-۱ الگوریتم پیشنهادی
۳۶	۴-۳-۲ تعیین توابع $f(\cdot)$ و $g(\cdot)$ در انتگرال فازی چوکت
۳۹	۴-۴ علت کارکرد انتگرال فازی چوکت در انتقال دانش
۴۱	۴-۴-۱ اثبات همگرایی روش پیشنهادی
۴۲	۵-۴ نتیجه‌گیری

#### فصل پنجم: نتایج شبیه‌سازی و آزمایش‌ها

۴۳	۱-۵ مقدمه
۴۴	۲-۵ رفتار الگوریتم‌های معرفی شده برای $g(\cdot)$
۴۵	۵-۲-۱ تعابیر مختلف انتگرال فازی چوکت از داده‌ها بر مبنای $g(\cdot)$
۴۶	۵-۳ مقایسه‌ی روش پیشنهادی با روش کوتاه‌ترین مسیر تجربه‌شده
۴۸	۵-۳-۱ مقایسه در محیط پلکان مارپیچ
۵۹	۵-۳-۲ مقایسه در محیط صید و صیاد
۶۷	۵-۴ بررسی تاثیر تعداد نواحی محیط در کیفیت و سرعت یادگیری عامل‌ها در روش پیشنهادی
۶۸	۵-۴-۱ محیط پلکان مارپیچ
۶۸	۵-۴-۲ محیط پلکان صید و صیاد
۶۸	۵-۵ بررسی تاثیر استفاده از انتگرال فازی در بهبود دانش جمعی
۷۲	۵-۶ تحلیل نتایج
۷۲	۵-۶-۱ مقایسه‌ی روش SEP با روش پیشنهادی

۷۴	۵-۶-۲ مقایسه‌ی تابع بولتزمن و $\varepsilon$ -حریصانه
۷۴	۵-۶-۳ بررسی تاثیر تعداد نواحی در کیفیت و سرعت یادگیری در روش پیشنهادی
۷۵	۵-۷ نتیجه‌گیری
۷۶	<b>فصل ششم: نتیجه‌گیری و جمع‌بندی</b>
۷۶	۶-۱ مقدمه
۷۷	۶-۲ نوآوری‌ها و نتایج کلی پایان‌نامه
۷۸	۶-۳ راهکارهای آینده و پیشنهادها
۷۸	<b>مراجع</b>
۸۱	<b>چکیده انگلیسی</b>

## فهرست تصاویر

- ۱-۱ جایگاه پژوهش انجام شده [۱، ۲] . . . . . ۳
- ۱-۲ شماتیک مکانیزم روش تخته‌سیاه برای یادگیری تقویتی مشارکتی [۲] . . . . . ۱۳
- ۲-۲ شمایی از یادگیری مشارکتی بر مبنای خبرگی عامل‌ها [۱] . . . . . ۱۵
- ۱-۳ شمایی از فرایند یادگیری تقویتی در تعامل با محیط [۲] . . . . . ۱۹
- ۲-۳ محیط پلکان مارپیچ [۲] . . . . . ۲۲
- ۳-۳ محیط صید و صیاد . . . . . ۲۳
- ۴-۳ دامنه‌ی دید و حالت تعریف شده برای عامل صیاد در محیط صید و صیاد [۲] . . . . . ۲۴
- ۵-۳ سرعت و کیفیت یادگیری از معیارهای ارزیابی و مقایسه‌ی عملکرد الگوریتم‌های یادگیری تقویتی می‌باشد [۲] . . . . . ۲۵
- ۱-۵ دو توزیع فرضی بجهت نمایش نحوه‌ی رفتار الگوریتم‌های ۴-۴ تا ۷-۴ بروی آن‌ها. . . . . ۴۴
- ۲-۵ نمایش توزیع‌های جدید بدست آمده بعد از اعمال الگوریتم‌های ۴-۴ تا ۷-۴ بروی دو توزیع فرضی شکل ۱-۵ . . . . . ۴۵
- ۳-۵ نمایش رفتار انتگرال فازی بروی منابع اطلاعاتی  $y = 1$  و  $y = 2$  و  $y = 3$  به ازای توابع  $g(\cdot)$ ‌های مختلف. . . . . ۴۶
- ۴-۵ مقایسه در سرعت و کیفیت یادگیری با تابع بولتزمن در محیط پلکان مارپیچ . . . . . ۴۸
- ۵-۵ مقایسه در سرعت اجرای روش‌ها به ازای تعداد تلاش‌های متفاوت برحسب میلی‌ثانیه با تابع بولتزمن در محیط پلکان مارپیچ . . . . . ۵۰
- ۶-۵ نمودار باروری الگوریتم‌ها مختلف با تابع بولتزمن در محیط پلکان مارپیچ . . . . . ۵۱
- ۷-۵ مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری با تابع بولتزمن در محیط پلکان مارپیچ . . . . . ۵۳
- ۸-۵ مقایسه در سرعت و کیفیت یادگیری با تابع  $\varepsilon$ -حریصانه در محیط پلکان مارپیچ . . . . . ۵۴
- ۹-۵ مقایسه در سرعت اجرای روش‌ها به ازای تعداد تلاش‌های متفاوت برحسب میلی‌ثانیه با تابع  $\varepsilon$ -حریصانه در محیط پلکان مارپیچ . . . . . ۵۵
- ۱۰-۵ نمودار باروری الگوریتم‌ها مختلف با تابع  $\varepsilon$ -حریصانه در محیط پلکان مارپیچ . . . . . ۵۶
- ۱۱-۵ مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری با تابع  $\varepsilon$ -حریصانه در محیط پلکان مارپیچ . . . . . ۵۷
- ۱۲-۵ مقایسه در سرعت و کیفیت یادگیری در محیط صید و صیاد با تابع بولتزمن در محیط صید و صیاد . . . . . ۵۹
- ۱۳-۵ مقایسه در سرعت اجرای روش‌ها به ازای تعداد تلاش‌های متفاوت برحسب میلی‌ثانیه با تابع بولتزمن در محیط صید و صیاد . . . . . ۶۱
- ۱۴-۵ نمودار باروری الگوریتم‌ها مختلف با تابع بولتزمن در محیط صید و صیاد . . . . . ۶۱

- ۱۵-۵ مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری با تابع بولتزمن در محیط صید و صیاد . . . . . ۶۲
- ۱۶-۵ مقایسه در سرعت و کیفیت یادگیری با تابع  $\varepsilon$ -حریصانه در محیط صید و صیاد . . . . . ۶۳
- ۱۷-۵ مقایسه در سرعت اجرای روش‌ها به ازای تعداد تلاش‌های متفاوت برحسب میلی‌ثانیه با تابع  $\varepsilon$ -حریصانه در محیط صید و صیاد . . . . . ۶۴
- ۱۸-۵ نمودار باروری الگوریتم‌ها مختلف با تابع  $\varepsilon$ -حریصانه در محیط صید و صیاد . . . . . ۶۵
- ۱۹-۵ مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری با تابع  $\varepsilon$ -حریصانه در محیط صید و صیاد . . . . . ۶۶
- ۲۰-۵ تاثیر ناحیه‌بندی مختلف بروی کیفیت و سرعت یادگیری در محیط پلکان مارپیچ . . . . . ۶۹
- ۲۱-۵ تاثیر ناحیه‌بندی مختلف بروی کیفیت و سرعت یادگیری در محیط صید و صیاد . . . . . ۷۰
- ۲۲-۵ تاثیر استفاده از انتگرال فازی در روش SEP بروی کیفیت و سرعت یادگیری در محیط پلکان مارپیچ. الف - استفاده از انتگرال فازی در روش SEP در محیط پلکان مارپیچ. ب - استفاده از انتگرال فازی در روش SEP در محیط صید و صیاد. . . . . ۷۱
- ۲۳-۵ تاثیر استفاده از انتگرال فازی در روش MCE بروی کیفیت و سرعت یادگیری در محیط پلکان مارپیچ. الف - استفاده از انتگرال فازی در روش MCE در محیط پلکان مارپیچ. ب - استفاده از انتگرال فازی در روش MCE در محیط صید و صیاد. . . . . ۷۳

## فهرست جداول

۲۰	۱-۳ ساختار جدول CP [۲]
۴۷	۱-۵ لیست اختصارهای استفاده شده در این فصل
۴۹	۲-۵ مقایسه در میزان درصد بهبود کیفیت یادگیری در محیط پلکان مارپیچ با تابع بولتزمن
۵۴	۳-۵ مقایسه در میزان درصد بهبود کیفیت یادگیری در محیط پلکان مارپیچ با تابع $\varepsilon$ -حریصانه
۵۷	۴-۵ مقایسه در سرعت و کیفیت یادگیری نسبت کیفیت نتیجه‌ی حاصل از تابع $\varepsilon$ -حریصانه نسبت به تابع بولتزمن
۵۸	۵-۵ مقایسه در نسبت میانگین سرعت اجرای حاصل از استفاده تابع $\varepsilon$ -حریصانه نسبت به تابع بولتزمن
۵۸	۶-۵ مقایسه در نسبت میزان باروری حاصل از استفاده تابع $\varepsilon$ -حریصانه نسبت به تابع بولتزمن
۵۹	۷-۵ مقایسه نسبت شیب تاثیر تعداد عامل‌ها میزان کیفیت نتیجه‌ی حاصل از تابع $\varepsilon$ -حریصانه نسبت به تابع بولتزمن
۶۰	۸-۵ مقایسه در میزان درصد بهبود کیفیت یادگیری در محیط صید و صیاد با تابع بولتزمن
۶۳	۹-۵ مقایسه در میزان درصد بهبود کیفیت یادگیری در محیط صید و صیاد با تابع $\varepsilon$ -حریصانه
۶۶	۱۰-۵ مقایسه در سرعت و کیفیت یادگیری نسبت کیفیت نتیجه‌ی حاصل از تابع $\varepsilon$ -حریصانه نسبت به تابع بولتزمن
۶۷	۱۱-۵ مقایسه در نسبت میانگین سرعت اجرا حاصل از استفاده تابع $\varepsilon$ -حریصانه نسبت به تابع بولتزمن
۶۷	۱۲-۵ مقایسه در نسبت میزان باروری حاصل از استفاده تابع $\varepsilon$ -حریصانه نسبت به تابع بولتزمن
۶۸	۱۳-۵ مقایسه در نسبت شیب تاثیر تعداد عامل‌ها میزان کیفیت نتیجه‌ی حاصل از تابع $\varepsilon$ -حریصانه نسبت به تابع بولتزمن

## چکیده

معمولا در دنیایی واقعی هنگامی که افراد برای انتقال دانش گرد هم می‌آیند و از تجربیات خوب و بد گذشته خود سخن می‌گویند هرکسی متناسب با جایگاهی که دارد دارای دانشی می‌باشد و در این انتقال دانش‌ها تجربیات هیچ کسی را نمی‌توان نادیده گرفت ولی گاهی پیش می‌آید که تجربیات و دانش فردی دارای بار محتویاتی بیشتری نسبت به اطرافیان خود می‌باشد، مردم معمولا از دانش فرد خبره‌تر بیشتر بهره می‌برند تا افراد دیگر. دستاوردهای این پژوهش بر مبنای همین فلسفه بنا شده است که سخن و دانش هرکسی باید شنیده شود. انتگرال فازی یکی از قوی‌ترین و منعطف‌ترین ابزارهای ریاضی برای ترکیب اطلاعات می‌باشد، لذا در این پژوهش از انتگرال فازی برای شنیدن بازتاب ندای دانش هر عامل در دانش جمعی و مدل کردن اطلاعات (دانش‌های) غیرافزایشی استفاده شده است. ولی در این راه مشکلاتی نیز وجود داشت و آن این بود که چگونه منصفانه بفهمیم که کدام عامل خبره‌تر از دیگری می‌باشد؟ در گذشته روش‌های متنوعی برای تخمین این معیار ارائه شد است که از شمارش میزان پاداش‌های مثبت و منفی عامل‌ها گرفته تا محاسبات پیچیده‌ای چون معیارهای شوک و کوتاه‌ترین مسیر تجربه شده. در طی پژوهش که منجر به نگارش این پایان‌نامه گردید احساس شد که تمامی روش‌های قبلی در یک چیز مشترکند: بسیار پیچیده و غیر منعطف!

وجود این فصل مشترک ناکارا انگیزه‌ای شد که در صدد ارائه‌ای معیاری برآیم که نه تنها ساده باشد بلکه در زندگی روزمره ما انسان‌ها هم تجلی داشته باشد. در پی این هدف ما به ارائه‌ی تئوری جامعی برای خبرگی پرداختیم که می‌تواند منشع بسیاری از تعاریف خبرگی، در آینده گردد؛ نهایتا با استفاده از تئوری خبرگی معرفی شده تعریفی برای یک معیار خبرگی جدید ارائه دادیم و نشان دادیم که تئوری و تعریف خبرگی جدید نسبت به تعاریف قبلی بسیار کارآمد بوده است.

**واژه‌های کلیدی:** ۱- سیستم‌های چندعامله، ۲- یادگیری مشارکتی، ۳- یادگیری تقویتی، ۴- دانش غیرافزایشی، ۵- انتگرال فازی.

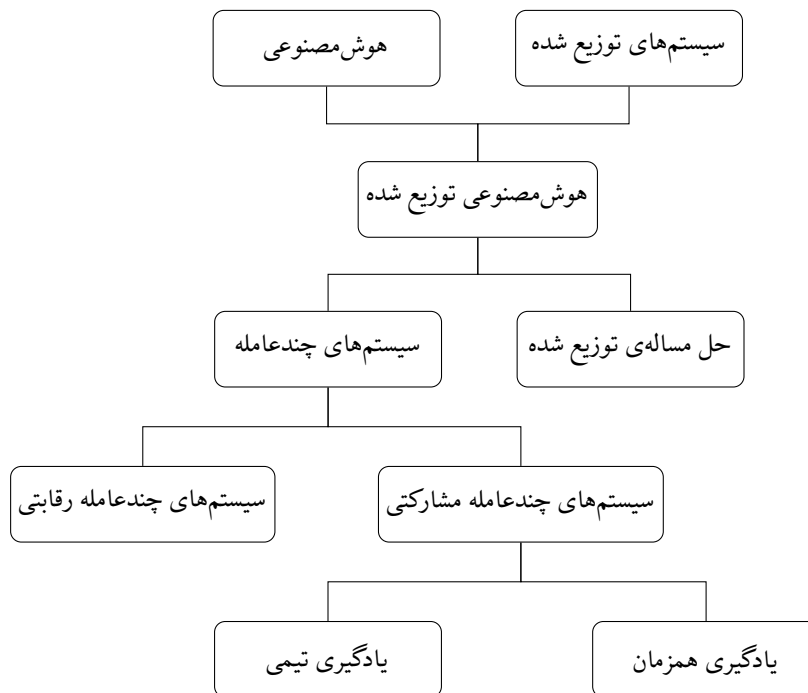
# فصل اول

## مقدمه

رایانه‌ها از قدرت بالایی در محاسبات برخوردار هستند تا جایی که محاسباتی که عامل انسانی ممکن است در چند سال انجام دهد را می‌توانند در کسری از ثانیه انجام دهند. در سال‌هایی که از عمر کامپیوتر گذشته پیشرفت‌های فراوانی صورت گرفته است که بر سرعت اجرای محاسبات کامپیوتر افزوده است. اما با وجود تمام این پیشرفت‌ها هنوز هم مسائل زیادی هستند که کامپیوتر نمی‌تواند در زمان قابل قبول آن‌ها را حل کند. یک روش که در راستای افزایش سرعت حل مسائل و با بهره‌گیری از زندگی جمعی عامل‌های انسانی پیشنهاد شد تقسیم مسئله به چندین زیر مسئله و حل هر یک توسط یک رایانه بوده است این موضوع شاخه‌ای تحت عنوان سیستم‌های توزیع شده را به وجود آورده است.

از طرف دیگر یکی از مسائل مهم در دنیای رایانه‌ها که از نظر محاسباتی به زمان زیاد نیاز دارد موضوعات هوش مصنوعی هستند که در ترکیب با سیستم‌های توزیع شده به هوش مصنوعی توزیع شده بدل شده‌اند. هوش مصنوعی توزیع شده نیز در قالب حل مسائل هوش مصنوعی و یادگیری فعالیت می‌کند. در حل مسائل توزیع شده یک مسئله که قابلیت اجرای موازی داشته باشد را به بخش‌هایی تقسیم کرده و هر بخش در رایانه‌ای حل شده و نهایتاً نتایج ترکیب می‌شوند این ترکیب می‌تواند یک یا چند مرتبه انجام شود.

دسته دیگر از مسائل هوش مصنوعی توزیع شده، مسائل یادگیری هستند که از زیر مجموعه‌ی سیستم‌های چندعاملی شناخته می‌شوند؛ این سیستم‌ها در راستای یادگیری از چند عامل بهره می‌برند. بر اساس روابطی که



شکل ۱-۱: جایگاه پژوهش انجام شده [۱، ۲]

می‌توان بین این عامل‌ها تعریف کرد سیستم‌های چندعاملی می‌توانند رقابتی یا مشارکتی باشند در سیستم‌های چند عامله رقابتی عامل‌ها به دنبال افزایش سود شخصی خود هستند که در زمان‌های زیادی به قیمت کاهش سود دیگر عامل‌ها خواهد بود؛ در سیستم‌های چندعاملی مشارکتی که موضوع پژوهش پیش رو نیز هست عامل‌ها به دنبال افزایش سود گروهی هستند. در شکل ۱-۱ جایگاه یادگیری مشارکتی در سیستم‌های چند عامله آورده شده است.

سیستم‌های چند عامله مشارکتی در دو دیدگاه مورد بررسی قرار می‌گیرند در دیدگاه اول عامل‌ها در یک محیط قرارگرفته و سعی در یادگیری نحوه‌ی تعامل با محیط و یکدیگر را دارند و این یادگیری در این راستا است که عامل‌ها بتوانند در محیط با همکاری هم به اهداف مشترک و مشخص شده برسند. در دیدگاه دوم عامل‌ها تعاملی ندارند یعنی که قرار نیست با همکاری یکدیگر به یک هدف مشخص دست‌یابند بلکه فقط سعی دارند با همکاری یکدیگر به یادگیری برسند این عامل‌ها در محیط‌های جداگانه و مشابهی قرار می‌گیرند و در طول فرایند یادگیری با هم ارتباط دارند. در این ارتباط داده‌های به دست آمده را به یکدیگر منتقل می‌نمایند تا زمانی که عامل‌ها تمام محیط را به‌خوبی مورد شناسایی قرار دهند.

#### ۱-۱ یادگیری مشارکتی در سیستم‌های چند عامله

انسان در طول حیات خود یادگیری زیادی انجام می‌دهد اما اگر قرار بود هر انسان از صفر شروع به جمع‌آوری اطلاعات کند و از عامل‌های دیگر یادگیری نداشته باشد بدون شک هنوز انسان‌ها همانند انسان‌های اولیه زندگی



می‌کردند. این رشدی که امروزه در زندگی انسانی دیده می‌شود مدیون انتقال اطلاعات و دانش بین عامل‌های انسانی است. بر همین اساس در [۳] یادگیری مشارکتی را سیستمی می‌داند که عامل‌ها در آن با همکاری یکدیگر به یادگیری یک وظیفه مشترک می‌پردازند. می‌توان آثار مثبت یادگیری مشارکتی در سیستم‌های چند عامله را چنین برشمرد.

- افزایش سرعت و دقت یادگیری.
- آزادسازی عامل از بهینگی محلی.
- کمک به تنظیم پارامترهای محلی عامل‌ها.

نکته مهمی که باید در نظر داشت عامل‌های موجود در یادگیری مشارکتی است. در بعضی از روش‌ها عامل‌های یک سیستم از توانایی‌های متفاوتی برخوردار هستند. به مجموعه عامل‌های موجود در این محیط‌ها تیم گفته می‌شود [۱، ۲] و در سیستم‌های دیگر مشابه پژوهش پیش رو عامل‌ها با توانایی‌های یکسان در نظر گرفته می‌شوند. گذشته از تفاوت بین عامل‌ها، یادگیری مشارکتی نیز همانند دیگر شاخه‌های هوش مصنوعی از چالش‌های فراوانی برخوردار است. فارغ از چالش‌های مشترکی که بیت یادگیری مشارکتی و روش یادگیری استفاده در آن وجود دارد چالش‌های جدیدی در یادگیری مشارکتی وجود دارد که می‌توان با سؤالاتی تعدادی از این چالش‌ها را نشان داد.

- چه اطلاعاتی باید بین عامل‌ها ردوبدل شود؟
- چه زمان باید اطلاعات منتقل شود؟
- ترکیب داده‌های دریافتی باید به چه صورت باشد؟
- بر اساس چه معیاری می‌توان عامل‌ها را مقایسه کرد؟

در رابطه با هریک از این چالش‌ها کارهای فراوانی چون پنددهی، تقلید و خبرگی انجام شده است که تعدادی از آن‌ها در فصل دوم آورده شده است. تمام این روش‌های تشریح شده در فصل دوم از یادگیری تقویتی به عنوان الگوریتم اصلی عامل برای یادگیری نحوه‌ی تعامل با محیط بهره برده‌اند و سعی در ارائه‌ی روش جهت ترکیب مناسب داده‌ها نموده‌اند. در ترکیب داده‌ها همیشه نیازی به معیاری جهت سنجش میزان درستی داده‌ها (دانش هر عامل) وجود دارد اما از آنجایی که عامل‌ها محیط را نمی‌شناسند در نتیجه دستیابی به معیاری صحیح برای این منظور کار دشواری می‌باشد.

در روش‌هایی چون خبرگی سعی شده تا معیارهایی جهت سنجش داده‌ها ارائه شود اما در کار پنددهی عامل‌ها زمانی که از داده‌ی خود مطمئن باشند به عامل دیگر بازخورد می‌دهند و این بازخورد در درک عامل از محیط موثر خواهد بود. معمولاً معیارهای معرفی شده در پژوهش‌های صورت گرفته بروی خبرگی یا با دیدگاه خیلی جزئی به بررسی خبرگی عامل‌ها می‌پردازند یا بصورت خیلی کلی؛ در حالت خلاصه کار اصلی که در این پژوهش انجام دادیم ارائه‌ی چهارچوب کلی برای تولید انواع معیارها و سپس ارائه‌ی معیاری که خیلی کلی یا جزئی نباشد و در عین حال بتواند عملکرد بهتری نسبت به روش‌های قبلی ارائه دهد.

## ۱-۲ اهداف و نوآوری‌های پایان‌نامه

بطور خلاصه هدف این پژوهش ارائه‌ی معیاری نرم (ساده در عین موثر بودن) بجهت محاسبه‌ی خبرگی عامل‌ها با در نظر گرفتن خاصیت غیرافزایشی [۴] خبرگی و دانش عامل‌ها می‌باشد - خاصیت غیرافزایشی دانش (خبرگی) عامل‌ها می‌گوید که ارزش دانش (خبرگی) چند عامل باهم لزوماً برابر با مجموع ارزش دانش (خبرگی) تک‌تک آن‌ها نمی‌باشد. در طی دستیابی به هدف تعیین شده در این پژوهش ابتدا چهارچوبی به نام «فرضیه‌ی خبرگی» معرفی شد که توانایی استخراج معیارهای زیادی برای محاسبه‌ی خبرگی از طریق این فرضیه میسر باشد؛ سپس با استفاده از فرضیه‌ی خبرگی معرفی شده و با در نظر داشتن هدف تعیین شده برای این پژوهش در مورد ارائه‌ی معیار نرم، معیار خبرگی جدیدی به نام «میزان ارجاع» تعریف شد.

در طی این پژوهش الگوریتمی برای ترکیب دانش عامل‌ها با در نظر داشتن میزان خبرگی معرفی شده هر عامل ارائه شد. در این الگوریتم از انتگرال فازی به عنوان عملگر ترکیب کننده دانش عامل‌ها استفاده کردیم و طبق آزمایش‌ها نشان دادیم که انتگرال فازی چوکت می‌تواند نتایج بهتری نسبت به روش‌های سنتی چون میانگیری وزن دار تولید کند زیرا انتگرال فازی چوکت می‌تواند خاصیت غیرافزایشی مساله را برخلاف میانگین وزنی مدل کند. دستاوردهای این پژوهش به صورت خلاصه به شرح زیر می‌باشد:

- معرفی چهارچوبی به نام «فرضیه‌ی خبرگی» برای تعریف معیارهای خبرگی جدید.
- تعریف معیار خبرگی جدید به نام «میزان ارجاع» در چهارچوب معرفی شده توسط «فرضیه‌ی خبرگی».
- استفاده از «انتگرال فازی چوکت» در ترکیب دانش‌های عامل‌ها با توجه به میزان خبرگی عامل‌ها.
- تعریف معیاری جدید به نام «میزان باروری» به جهت سنجش سرعت یادگیری الگوریتم‌ها.
- بررسی تاثیر سیاست‌های انتخاب عمل  $\epsilon$  - حریصانه<sup>۱</sup> در یادگیری مشارکتی - پژوهش‌های قبلی این موضوع را مورد بررسی قرار نداده‌اند.

<sup>۱</sup>  $\epsilon$ -greedy

- اثبات صحت فرضیه و معیار خبرگی معرفی شده در این پژوهش با توجه نتایج آزمایش‌ها.

### ۳-۱ ساختار پایان‌نامه

در ادامه‌ی گزارش در فصل دوم سعی شده کارهای انجام‌شده در این زمینه تشریح شود؛ در فصل سوم موضوعاتی که برای درک روش پیشنهادی در این پژوهش لازم است بیان شده است؛ سپس در فصل چهارم روش پیشنهادی تشریح و در فصل پنجم آزمایش‌های موردنیاز جهت نمایش عملکرد روش پیشنهادی آورده شده است. نهایتاً در فصل ششم جمع‌بندی از مطالب ارائه شده در این پایان‌نامه صورت گرفته است.

## فصل دوم

### مرور کارهای پیشین

#### ۲-۱ مقدمه

در سال‌های گذشته پژوهش‌های فراوانی در سیستم‌های چند عامله انجام شده است که در این پژوهش‌ها محققان سعی داشته‌اند مزایای کار گروهی در انسان را در رایانه نیز ایجاد نمایند. یکی از قابلیت‌های عامل‌های هوشمند که می‌تواند با کار گروه سریع‌تر و بهتر شود موضوع یادگیری است که در این زمینه هم کارهایی انجام شده که معمولاً الگوبرداری از عامل‌های انسانی بوده است. همان‌طور که می‌دانیم انسان تنها از یک مکانیزم در رسیدن به یادگیری بهره نمی‌برد، عامل‌های انسانی با تقلید از عامل‌هایی که دارای اطلاعات بیشتری هستند توانسته‌اند یادگیری خود را بهبود دهند؛ عامل‌های انسانی در شرایط بحرانی زندگی از عامل‌های باتجربه‌تر پند می‌گیرند، عامل‌های انسانی در مراتبی از خبرگی قرار دارند؛ همه این موارد الگوهای مناسب بوده که توانسته یادگیری در سیستم‌های چندعاملی را بهبود بخشد. اما می‌توان کارهایی که در یادگیری مشارکتی انجام می‌شود را به دسته‌هایی تقسیم کرد، هر دسته از پژوهش‌های انجام شده در این رشته سعی در رفع یک یا چند چالش از چالش‌های این رشته داشته‌اند. پژوهش پیش رو را می‌توان از دسته پژوهش‌های یادگیری مشارکتی دانست که سعی در حل مشکل ترکیب داده‌های عامل‌ها دارند که روش‌های ارائه شده در این فصل نیز روش‌هایی هستند که در تقسیم داده‌های یادگیری مشارکتی فعالیت کرده‌اند. پیچیدگی ترکیب داده‌های عامل به این دلیل است که معیار مناسبی جهت مشخص کردن داده‌ی درست وجود ندارد. در بسیاری از کارهایی که در این فصل ارائه خواهد شد در ترکیب

داده‌ها معیار جایگزینی معرفی شده و آن معیاری جهت نمایش برتری عامل است. ایده این جایگذاری از آنجاست که عاملی که از برتری برخوردار باشد داده‌های بهتری نیز نسبت به عامل‌های دیگر خواهد داشت.

## ۲-۲ اشتراک‌گذاری اطلاعات

برای اولین بار در [۵] اشتراک‌گذاری داده‌ها در سیستم‌های چند عامله مورد ارزیابی قرار گرفت. هدف این بررسی نمایش اثر اشتراک‌گذاری داده‌ها در مقابل سیستم‌های تک عاملی بود. نتیجه این پژوهش نشان داد که اگر اشتراک‌گذاری به خوبی انجام شود می‌تواند سرعت و کیفیت یادگیری را به صورت چشم‌گیری افزایش دهد. در این پژوهش سه نوع اشتراک‌گذاری مورد بررسی قرار گرفت در نوع اول که اشتراک‌گذاری ادراک نام گرفت عامل‌ها تنها نتایج مشاهدات خود را به اشتراک می‌گذاشتند، در نوع دوم اشتراک‌گذاری سه‌تایی حالت، عمل، کیفیت اشتراک‌گذاری شده و اشتراک‌گذاری واقعیت نامیده شد و نهایتاً در نوع سوم اشتراک‌گذاری که اشتراک‌گذاری سیاست خوانده می‌شود اطلاعات داخلی عامل‌ها که منبع استخراج سیاست آنهاست به اشتراک گذاشته شده است. اشتراک‌گذاری در این پژوهش با یک میانگین‌گیری ساده بین اطلاعات عامل‌ها انجام می‌شد. در این پژوهش که SA<sup>۱</sup> نامیده شد ثابت شده ممکن است اشتراک‌گذاری سربارهایی در ترکیب داده‌ها به سیستم بیفزاید یا در شروع یادگیری از سرعت یادگیری بکاهد اما در طول یادگیری این سربارها جبران شده و اشتراک داده‌ها می‌تواند به صورت چشم‌گیری در افزایش سرعت سیستم‌های چند عامله مؤثر باشد.

## ۳-۲ یادگیری مشترک

برنجی و همکاران در سال ۱۳۷۸ (۱۹۹۹ م.) روشی تحت عنوان یادگیری مشترک مطرح کردند [۶]. در این روش اشتراک‌گذاری با در نظر گرفتن تنها یک سیاست برای تمام عامل‌ها انجام شد. نتایج این پژوهش نشان می‌دهد که در دسته بزرگی از مسائل روش‌های یادگیری مشترک می‌تواند مفیدتر از روش‌های یادگیری مستقل باشد. فرآیند یادگیری در این روش به این صورت است که عامل‌ها در محیط اقداماتی انجام می‌دهند و بعد از دریافت پاداش عمل بروزرسانی را در یک داده مشترک انجام می‌دهند و در انتخاب عمل نیز از همان داده مشترک بهره می‌برند. این به این معنی است که عامل‌ها دیگر برای خود داده مستقلی ندارند. در این پژوهش حتی یادگیری با منطق فازی ادغام شده است و نویسندگان سعی کردند اثر فازی کردن داده‌ها در یادگیری مشارکتی را نمایش دهند.

<sup>۱</sup> Simple Averaging

## ۲-۴ تقلید

انسان در طول زندگی برای رسیدن به یادگیری روش‌های متفاوتی دارد. گاهی برای رسیدن به یادگیری باید آزمایش کرد گاهی تحلیل کرد و گاهی تجربه اما یک روش که انسان از آن مخصوصاً در مراحل رشد بسیار بهره می‌برد تقلید است. همین موضوع باعث شده که در یادگیری مشارکتی نیز به تقلید عامل‌ها از هم توجه شود. بر همین اساس نونس و همکاران با ایده برداری از تقلید در انسان پیشنهاد کردند که رابطه عامل‌ها از طریق تقلید از یکدیگر باشد [۷].

موضوع دیگری که در مورد تقلید عامل‌های انسانی باید در نظر گرفته می‌شد این است که عامل‌های انسانی از عامل‌های انسانی تقلید می‌کنند که اطلاعات بیشتری دارند. در پیاده‌سازی انجام‌شده نیز بر همین اساس سه نوع تقلید پیشنهاد می‌شود. تقلید می‌تواند به صورت ساده باشد. پیشنهاد داده‌شده است که عامل‌ها همیشه از عامل‌های همسایه (همسایگی در این روش بر اساس همسایگی محلی است چرا که عامل‌هایی که در منطقه یکسانی قرار دارند کمک بیشتری می‌توانند به هم کنند) خود تقلید نمایند. این موضوع یک دور در عامل‌ها ایجاد می‌کند که هر عامل منتظر می‌ماند تا عامل دیگر حرکتی انجام دهد. برای رفع این موضوع نوع دیگری از تقلید به نام تقلید شرطی مطرح می‌شود در تقلید شرطی عامل از کسانی تقلید می‌کند که عملکرد بهتری نسبت به او داشته‌اند در این حالت موضوع دور و انتظار عامل‌ها برطرف شده است. اما در روش سوم که تقلید انطباقی نام دارد عامل همیشه تقلید نکرده و تقلید بر اساس یک احتمال انجام می‌شود.

## ۲-۵ حافظه جمعی

گارلند و همکاران در سال ۱۳۷۵ (۱۹۹۶ م.) ایده جدید خود را با عنوان یادگیری حافظه جمعی مطرح کردند [۸، ۹]. در یادگیری حافظه جمعی که برگرفته از شناخت توزیع‌شده در علوم اجتماعی می‌باشد عامل‌ها تجارب خود را در یک حافظه مشترک نگهداری می‌کنند. هر عامل در زمان برخورد با مشکلات می‌تواند با بهره‌گیری از این تجارب راه درست را پیدا کند. این روش در دو دیدگاه مورد ارزیابی قرار گرفته است. در دیدگاه اول عامل‌ها الگوهای موفق خود در طول یادگیری را در حافظه مشترک نگهداری می‌کنند تا در زمان نیاز تمام عامل‌ها با استفاده از این الگوها بتوانند راه‌حل مشکلات خود را پیدا کنند. در دیدگاه دیگر احتمال موفقیت عامل‌ها نگهداری می‌شود که با بهره‌گیری از این داده می‌توان میزان موفقیت عامل‌ها در اعمال مختلف را ارزیابی کرده و در جهت بهبود طراحی سیستم مورد ارزیابی قرارداد. لازم به ذکر است که در این پژوهش‌ها حافظه جمعی را در دو حالت حافظه مرکزی و حافظه توزیع‌شده بین عامل‌ها مورد ارزیابی قرار داده است.

## ۲-۶ پند

در سال ۱۳۸۱ (۲۰۰۲ م.) نونس و همکاران باردیگر روشی جدید با عنوان پند دهی مطرح کردند [۱۰]. در جوامع انسانی پند دادن بسیار رواج داشته و در زمان مشکلات بسیار کارا می‌باشد. یک عامل انسانی در زمان برخورد با مشکلات از عامل‌هایی که اطلاعات بیشتری دارند پند گرفته و مشکلات خود را حل می‌کند. عاملی انسانی که دارای اطلاعاتی است هم اطلاعات خود را با تجربه کردن و یا گرفتن پند در زمان‌های دیگر به‌دست می‌آورد. مشخصاً یادگیری تقویتی در حالت معمول با تجارب به یادگیری می‌رسد. اگر هر تجربه را بازخوردی از محیط در نظر بگیریم هر پند را نیز می‌توان بازخوردی از عامل‌های دیگر دانست. با این ایده دیگر حتی نیازی نیست که عامل‌ها از روش‌های یکسانی در یادگیری بهره ببرند زیرا پند دادن به عامل‌ها را می‌توان فارغ از روش یادگیری پیاده‌سازی کرد. ایده پردازان پند در [۱۱] کار قبل خود را کامل‌تر کرده و این ایده را به‌صورتی که عامل‌ها در یک محیط به تعامل می‌پردازند پیاده‌سازی کردند. هر عامل بعد از رسیدن به هر حالت موقعیت خود را به عامل‌های دیگر ارسال می‌نماید. عامل‌هایی که تجربه مشابهی داشته‌اند در پاسخ مقداری را به عنوان میزان ارزش عمل انجام شده برای عامل ارسال می‌کنند و عامل از این مقادیر همانند پاداش دریافتی از محیط بهره می‌برد.

## ۲-۷ یادگیری مشارکتی بر مبنای خبرگی

تشریح یادگیری مشارکتی بر مبنای خبرگی را با یک سؤال می‌توان آغاز کرد. آیا عامل‌ها در شناخت محیط از خبرگی یکسانی برخوردار هستند؟ مسلماً چنین نیست، در [۱۲] ایده یادگیری مشارکتی بر مبنای خبرگی با عنوان WSS<sup>۱</sup> مطرح می‌شود. همان‌طور که در تشریح روش SA مطرح شد در این روش با میانگین‌گیری از اطلاعات عامل‌ها ترکیب انجام می‌شود. در این میانگین‌گیری تمام عامل‌ها به یک اندازه سهم هستند. ایده پردازان WSS با طرح این موضوع که میزان خبرگی عامل‌ها یکسان نیست سعی کردند هر عامل در ترکیب داده‌ها به میزان توانایی و خبرگی خودش مؤثر باشد.

نویسندگان با ارائه معیارهایی میزان خبرگی عامل‌ها را سنجیده و بر همین اساس داده‌ها باهم ترکیب می‌شوند. در WSS روال یادگیری به دو فاز یادگیری مستقل و یادگیری مشارکتی شکسته شده است. در یادگیری مستقل هر عامل به‌طور مستقل به یادگیری می‌پردازد این یادگیری منجر به کسب اطلاعاتی می‌شود که در فاز یادگیری مشارکتی باهم ترکیب می‌شوند. یادگیر در فاز یادگیری مستقل چندین چرخه یادگیری را تجربه می‌کند. تعداد این چرخه‌ها می‌تواند در بین عامل‌ها یکسان و یا متفاوت باشد. اما باید در انتخاب تعداد چرخه‌های یادگیری هر فاز یادگیری مستقل دقت کرد چراکه اگر این تعداد کم در نظر گرفته شود عامل اطلاعات کافی را جمع‌آوری

<sup>۱</sup> Weighted Strategy Sharing

نکرده است و اگر زیاد در نظر گرفته شود از تأثیر یادگیری مشارکتی خواهد کاست.

در فاز دوم یادگیری عامل‌ها باید به یادگیری مشارکتی بپردازند. در آغاز این فاز میزان خبرگی عامل‌ها سنجیده می‌شود و پس از آن داده‌ها ترکیب‌شده و جداول  $Q$  عامل‌ها بروز رسانی می‌شود. در [۱۲] روش‌هایی جهت ترکیب داده‌ها ارائه شده است. در یکی از روش‌ها جدول تمام عامل‌ها با بهره‌گیری از میزان خبرگی میانگین‌گیری شده و جدول تولیدشده به تمام عامل‌ها داده شود که در صورت انجام این کار بعد از فاز یادگیری مشارکتی تمام عامل‌ها جدول  $Q$  یکسانی خواهند داشت. در روش دیگری پیشنهادشده که هر عامل جدول جدید خود را با ترکیب جدول خود با جدول عامل‌های خبره‌تر از خودش تولید کند. در این ترکیب نیز هر عامل به میزان خبرگی خودش در ترکیب داده‌ها سهم خواهد داشت.

در WSS با در نظر گرفتن خبرگی عامل‌ها تأثیر زیادی در بهبود یادگیری مشارکتی داشته است اما نکته‌ای که در نظر گرفته نشده است اینجاست که میزان خبرگی عامل‌ها در دامنه‌های مختلف بسیار متفاوت بوده و بهتر است که در ترکیب داده‌ها این دامنه‌ها هم در نظر گرفته شود. در [۱۳] با در نظر گرفتن دامنه خبرگی عامل‌ها سعی شده تا نقصان WSS برطرف شود. بعد از آن در [۱۴] سعی شده تا استفاده از جدول  $Q$  یک عامل در ترکیب داده‌ها قطعی نباشد. در این راستا در فاز ترکیب برای اطلاعات هر عامل احتمالی در نظر گرفته شده است که نشان‌دهنده احتمال حضور اطلاعات آن عامل در ترکیب داده‌ها است. میزان این احتمال نیز بر اساس تفاوت میزان خبرگی عامل‌ها محاسبه شده است. در ادامه تعدادی از معیارهای خبرگی معرفی شده در [۱۲] خواهد آمد.

• **معیار خبرگی معمولی:** در این معیار میزان خبرگی عامل‌ها بر اساس مجموع پاداش‌های دریافتی آنها در نظر گرفته شده است. در نتیجه عاملی که میزان پاداش منفی کمتر و میزان پاداش مثبت بیشتری گرفته است را عامل خبره‌تر می‌داند.

• **معیار خبرگی مثبت:** در این معیار سعی شده با شمارش پاداش‌های مثبت عامل‌ها میزان خبرگی اندازه‌گیری شود. ایده انتخاب این معیار این بوده که عاملی که پاداش مثبت بیشتری گرفته است از خبرگی بالاتری برخوردار است.

• **معیار خبرگی منفی:** این معیار برعکس معیار خبرگی مثبت با این ایده که عاملی که پاداش منفی بیشتری دارد نقاط بحرانی بیشتری را می‌شناسد عمل شده و تعداد پاداش‌های منفی عامل یادگیری را شمارش می‌نماید.

• **معیار خبرگی قدر مطلق:** در معیار خبرگی قدر مطلق میزان خبرگی عامل با محاسبه مجموع قدر مطلق پاداش‌ها دریافتی او انجام می‌شود. در نتیجه به پاداش‌های منفی و مثبت ارزش یکسانی داده شده است.



• **معیار خبرگی گرادیان:** در این معیار مانده معیار اول عمل می‌شود با این تفاوت که میزان افزایش سیگنال دریافتی نسبت به آخرین دوره‌ی یادگیری مشارکتی را معیار خبرگی قرار داده است، هرچقدر این اختلاف بیشتر و مثبت باشد نشان می‌دهد که عامل نسبت به دوره‌ی قبل خبره‌تر شده است.

• **معیار خبرگی میانگین تعداد قدم‌ها:** این معیار برعکس پنج معیار دیگر به‌جای تأکید بر روی پاداش‌ها میانگین تعداد قدم‌های عامل در چرخه‌های یادگیری را معیار می‌داند. این انتخاب با این ایده انجام‌شده که عامل‌های خبره‌تر با تعداد قدم‌های کمتر چرخه‌های یادگیری را به اتمام می‌رسانند.

## ۸-۲ یادگیری مشارکتی بر مبنای تخته‌سیاه

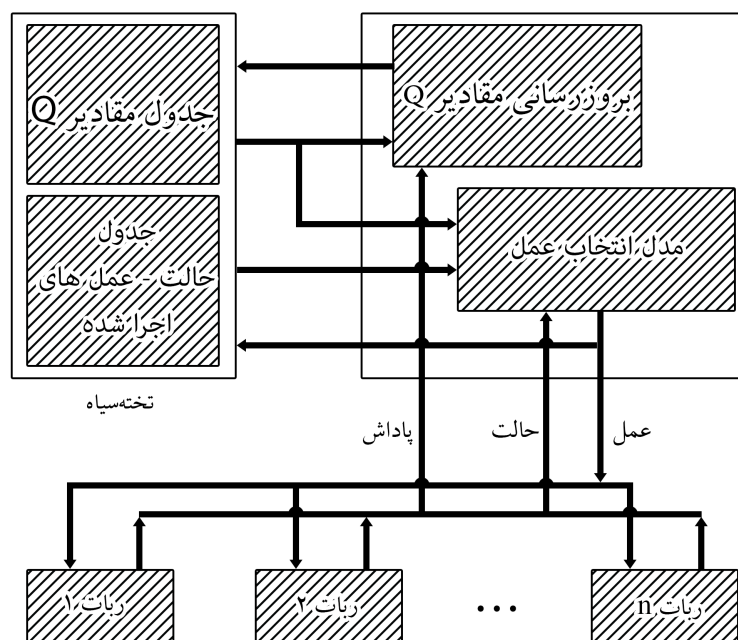
در [۱۵] سازوکار تخته‌سیاه مطرح شد. تخته‌سیاه یک حافظه مرکزی است که تمام عامل‌ها به آن دسترسی دارند. در این روش عامل‌ها به‌طور مستقیم باهم ارتباط نداشته و ارتباطات از طریق همین تخته‌سیاه انجام می‌شود. هر عامل می‌تواند بر روی تخته نوشته و یا از آن بخواند. در روش پیشنهادشده در [۱۵] به این شکل است که عامل بعد از رسیدن به هر موقعیت حالت خود را به تخته‌سیاه اعلام می‌کند و تخته‌سیاه عملی را بر اساس حالت جاری به عامل برمی‌گرداند. عامل بعد از انجام آن عمل و دریافت بازخورد از محیط این بازخورد را به تخته‌سیاه برمی‌گرداند.

تخته‌سیاه دودسته از داده‌ها را نگهداری می‌کند. دسته اول داده‌ها همان جدول  $Q$  عامل‌ها است و دسته دوم از داده‌ها عمل‌های انجام‌شده توسط هر عامل است. همان‌طور که مشخص است در این روش بروز رسانی جدول  $Q$  و انتخاب عمل از عامل به تخته‌سیاه منتقل شده و مشخصاً جدول  $Q$  باید در تخته‌سیاه پردازش شود. اما دسته دوم اطلاعات صرفاً جهت کمک به انتخاب عمل عامل‌ها انجام می‌شود. به‌عنوان مثال اگر عامل در حالتی قرار گیرد و عملی تجربه نشده باشد آن عمل پیشنهاد می‌شود. پس ذخیره‌سازی دسته دوم اطلاعات در جهت مدیریت اکتشاف و بهره‌برداری عامل‌ها از اطلاعات است. در شکل ۲-۱ مکانیسم تخته‌سیاه نمایش داده‌شده است.

## ۹-۲ یادگیری تقویتی تعاملی

در سال ۱۳۸۵ (۲۰۰۶ م.) لیما و همکاران سه الگوریتم به نام‌های  $Q$ -AVG,  $Q$ -BEST و  $Q$ -PSO را ارائه دادند که این دسته از الگوریتم‌ها مسأله‌ی یادگیری مشارکتی تقویتی را به صورت جستجوی جدول  $Q$  توسط عامل‌ها تعریف کرده است [۱۶]. لیما و همکاران ابتدا معیاری ارائه دادند برای سنجش میزان خوب بودن مقادیر جداول  $Q$  که در واقع متشکل از مقادیر پاداش‌های دریافتی عامل‌ها می‌باشد، که این معیار شامل جمع تخفیف یافته<sup>۱</sup> پاداش‌ها در هر چرخه‌ی یادگیری می‌باشد به‌صورتی که به پاداش‌های نهایی ارزش بیشتری می‌دهد. سپس با

<sup>۱</sup> Discounted Sum



شکل ۲-۱: شماتیک مکانیزم روش تخته‌سیاه برای یادگیری تقویتی مشارکتی [۲]

استفاده از این معیار به جستجوی بهترین جدول  $Q$  از طریق عامل‌ها می‌پردازد؛ که این مساله وجه تمایز کار لیما با روش پیشنهادی در این پژوهش می‌باشد. الگوریتم  $BEST-Q$  در مرحله‌ی به اشتراک‌گذاری دانش، بهترین دانش (با توجه به معیار معرفی شده) را به عنوان دانش جمعی در نظر می‌گیرد. الگوریتم  $AVG-Q$  بهترین دانش را با دانش هر عامل میانگین‌گیری می‌کند و به عنوان دانش همان عامل در نظر می‌گیرد و در الگوریتم  $PSO-Q$  الگوریتم  $PSO$  [۱۷] را به عنوان مدل جستجو کننده دانش جمعی در نظر گرفته است.

یکی از معایب این روش‌ها این است که یادگیری تقویتی را به صورت یک مساله‌ی جستجوی مقادیر جداول  $Q$  در نظر گرفته است و از آنجایی که ماهیت الگوریتم‌ها به صورت جستجو می‌باشد، اثباتی جهت اینکه این جستجوهای مقادیر جداول  $Q$  به مقادیر بهینه یعنی جدول  $Q^*$  همگرا خواهند شد، وجود ندارد. از دیگر معایب این روش‌ها نحوه‌ی محاسبه‌ی میزان بهینگی مقادیر جداول عامل‌ها می‌باشد، بطوری که این روش‌ها به پاداش‌های دریافتی در اواخر چرخه‌ی یادگیری عامل‌ها ارزش بیشتری می‌دهند، این ممکن است در ابتدای یادگیری که عامل‌ها در اوایل چرخه‌ی یادگیری خود حرکت‌های بی‌په‌وده‌ی زیادی انجام دهند منطقی به نظر بیاید، ولی بعد از آنکه تعدادی چرخه‌ی یادگیری سپری شد و عامل به دانش نسبی خوبی از محیط خود دست یافت حرکت‌های ابتدایی به اندازه‌ی حرکت‌های نهایی ارزش دارند (و البته شاید هم ارزش بیشتری داشته باشند) زیرا که عامل برای دست یافتن به هدف، نسبت به محیط اطراف اهداف شناخت بهتری دارد (به علت خاصیت شوک یادگیری تقویتی [۲]) که در این شرایط در حالت کلی میزان بهینگی عامل‌ها را بیشتر میزان بهینگی اعمال در اوایل چرخه‌ی یادگیری عامل تعیین می‌کند. معیاری که روش‌های فوق‌الذکر از آن، جهت سنجش خبرگی عامل‌ها در

نظر گرفته‌اند این مساله را نادیده می‌گیرد که باعث می‌شود خبرگی عامل‌ها را نتوان به درستی تعیین کرد.

## ۲-۱۰ یادگیری مشارکتی بر مبنای پختگی سیاست

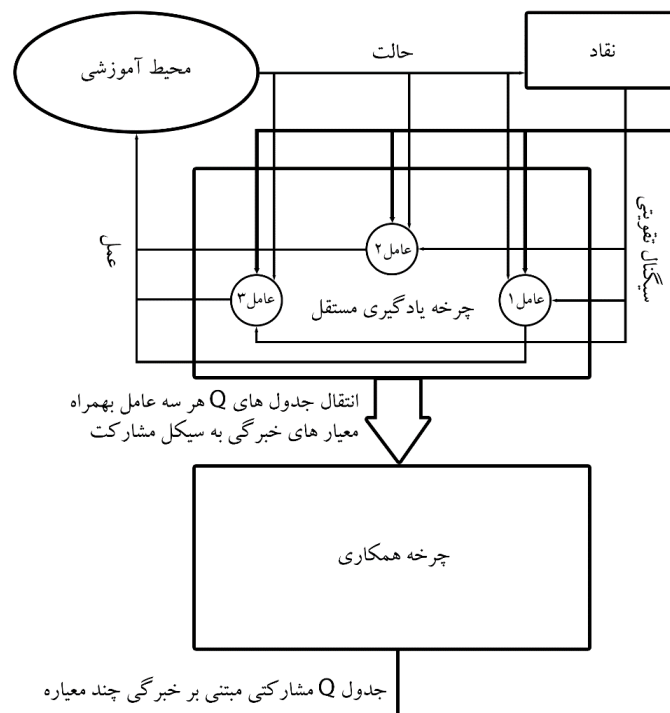
یانگ و همکاران در سال ۱۳۸۸ (۲۰۰۹ م.) روشی با عنوان یادگیری مشارکتی بر مبنای خبرگی چند معیاری ارائه دادند [۱۸]. این روش تا حدودی ترکیب روش تخته‌سیاه با WSS هست. در این روش عامل‌ها حافظه مرکزی خود یا تخته‌سیاه را دارند که وجود تخته‌سیاه عامل‌ها را از شکستن بازه یادگیری به دو فاز بی‌نیاز می‌سازد. در روشی چون WSS یادگیری به دو فاز یادگیری مستقل و یادگیری مشارکتی شکسته می‌شد تا عامل‌ها داده‌های خود را به اشتراک بگذارند اما زمانی که عامل‌ها دائماً می‌توانند داده‌های خود را بر روی تخته‌سیاه نوشته و بخوانند ارتباط از طریق همین تخته‌سیاه انجام خواهد شد.

اما عامل‌ها بر خلاف روش تخته‌سیاه از ارتباط مستقیم هم در تصمیم‌گیری‌ها و انتخاب اعمال بهره می‌برند. عاملی که در وضعیت انتخاب عمل قرار گرفته می‌تواند از عامل‌های دیگر بیاموزد. در این روش جهت شناخت عامل‌هایی که اطلاعات خوبی دارند و می‌توانند آموزگار باشند از معیارهای خبرگی ارائه‌شده در WSS استفاده شده است. با این کار عامل از عامل‌هایی می‌آموزد که واقعاً از خبرگی بالاتری برخوردار هستند. این کار باعث می‌شود که در شروع یادگیری که عامل‌ها داده کمی دارند و نیز عامل آموزگاری نداشته با کمک اطلاعات و دستورات تخته‌سیاه عمل کند و بعد طی مراحل یادگیری که عامل‌ها داده‌های زیادی کسب کردند با بهره بردن از نظرات دیگر عامل‌ها انتخاب‌های بهتری داشته باشند.

## ۲-۱۱ یادگیری مشارکتی بر مبنای خبرگی چند معیاری

پاکیزه و همکاران در سال ۱۳۹۲ (۲۰۱۳ م.) با نقد روش WSS روشی جدید ارائه کردند [۱۹]. ایشان با اشاره به این موضوع که خبرگی در یک‌رشته نبوده در کار خود از ترکیب ۶ معیار خبرگی WSS در کنار هم بهره برده‌اند. ایشان تأکید دارند که عامل‌های انسانی در زمینه‌های مختلف خبرگی‌های متفاوتی دارند و این موضوع در عامل‌های هوشمند نیز وجود دارد. این پژوهش هر یک از معیارهای ارائه‌شده در WSS را مانند یک زمینه در عامل انسانی دانسته و در روش خود از تمام این معیارها در کنار هم بهره برده‌اند.

پاکیزه و همکاران مانند WSS یادگیری را در دو فاز یادگیری مستقل و یادگیری مشارکتی تقسیم می‌نمایند عامل‌ها در فاز یادگیری مشترک از هر معیار برای ترکیب داده‌های جدول Q بهره می‌برند و بعد از ترکیب جدول به وسیله هر معیار ۶ جدول مشارکتی تولید می‌شود که هر یک بر اساس یک معیار خبرگی است. آن‌ها برای ترکیب این جداول آن‌ها را باهم جمع می‌کنند. اما موضوعی که وجود دارد این است که جدول تولیدشده به وسیله جمع چندین جدول دیگر خواص جدول Q را ندارد. برای رفع این مشکل این جدول را نه در جایگزینی با جدول Q



شکل ۲-۲: شمایی از یادگیری مشارکتی بر مبنای خبرگی عامل ها [۱]

عامل ها بلکه در کنار جدول  $Q$  عامل نگه داری می نمایند. به عبارت دیگر هر عامل دو جدول دارد یک جدول  $Q$  که بر اساس یادگیری تقویتی است و جدول دیگر که جدول مشارکتی عامل ها است. پاکیزه و همکاران پیشنهاد کردند که از جدول مشارکتی که خواص جدول  $Q$  عامل ها را ندارد صرفاً برای انتخاب عمل استفاده شود و عامل بر اساس این جدول عمل را انتخاب کرده انجام دهد سپس جدول  $Q$  خود را بروز رسانی نماید. جهت درک بهتر این روش شمای کلی آن در شکل ۲-۲ آورده شده است.

## ۱۲-۲ تسریع یادگیری مشارکتی با بهره گیری از کوتاه ترین فاصله تجربه شده

میرزایی در سال ۱۳۹۵ (۲۰۱۶ م.) جهت تسریع در یادگیری مشارکتی دو معیار جدید را ارائه کرد [۲]. معیار اول یک معیار مکاشفه است که کوتاه ترین فاصله تجربه شده توسط عامل از هر حالت و عمل را شمارش می کند. ایشان نام این معیار را SEP<sup>۱</sup> گذاشته است. معیار دیگر که شوک نام گذاری شده است میزان شناخت عامل از هر حالت و عمل را محاسبه می نماید.

میرزایی برخلاف دیگران فقط در فاز ترکیب داده های یادگیری مشارکتی تغییر ایجاد نکرده است. وی در فاز انتخاب عمل توسط عامل های مشارکتی نیز از جدول SEP در کنار جدول  $Q$  استفاده کرده است. استدلال ایشان در انجام این کار چنین بوده که عامل های یادگیری تقویتی در فازهای اول یادگیری داده زیادی ندارند و از آنجایی که جدول SEP با سرعت بیشتری به روز رسانی می شود بهتر است انتخاب اعمال در فازهای اولیه یادگیری

<sup>۱</sup> Shortest Experienced Path

بیشتر بر اساس SEP انجام شود. ایشان با استفاده از شوک که نمایشی از میزان شناخت عامل از هر حالت و عمل است تعادلی بین بهره‌برداری از جدول SEP و جدول  $Q$  برقرار کرده است. در شروع یادگیری که شناخت عامل کمتر است بیشتر انتخاب بر اساس SEP انجام می‌شود و در طول یادگیری با افزایش میزان شناخت عامل از محیط انتخاب عمل بر اساس جدول  $Q$  افزایش می‌یابد.

همچنین میرزایی در پژوهش خود در فاز ترکیب داده‌ها نیز روش جدیدی ارائه داد. از آنجایی که وی یک جدول جدید به سیستم افزوده است در فاز ترکیب داده‌ها جدول SEP عامل‌ها را نیز ترکیب می‌نماید. همچنین جداول SEP عامل‌ها را تنها با یک حداقل‌گیری باهم ترکیب کرده و به عامل‌ها برمی‌گرداند. سپس ترکیب جداول  $Q$  عامل‌ها به صورت محلی انجام می‌گیرد به این صورت که هر سطر از جدول که نمایش یک حالت از محیط است به صورت جداگانه بروز رسانی می‌شود. ایشان در ترکیب داده‌های هر سطر عامل‌ها را به دو گروه تقسیم نموده و داده‌های هر گروه را جداگانه ترکیب می‌نماید. این تقسیم‌بندی بر اساس رابطه بین سیاست‌های استخراج‌شده از جدول  $Q$  و SEP عامل در یک حالت هست. وی عامل‌هایی که سیاست استخراج‌شده از جدول  $Q$  و SEP در آنها همخوانی داشته باشد در یک گروه و عامل‌هایی که سیاست استخراج‌شده آنها عمل‌های متفاوتی را پیشنهاد می‌کنند را در گروه دیگر قرار داده است. ترکیب داده‌های هر گروه با استفاده از میزان شناخت عامل از آن حالت (شوک) انجام می‌شود به این صورت که داده‌های عملی که شناخت بیشتری دارند بیشتر مورد استفاده قرار می‌گیرند. در فصل بعد روش محاسبه ارائه شده توسط میرزایی تشریح شده است.

## ۱۳-۲ نتیجه‌گیری

روش‌های ارائه‌شده در این فصل در ترکیب داده‌ها در یادگیری مشارکتی عمل می‌نمایند. در بعضی از این روش‌ها سعی شده رابطه عامل‌ها به صورت غیرمستقیم باشد و در روش‌هایی چون WSS سعی شده تا زمانی برای جمع‌آوری داده‌ها به عامل داده شود و بعد از آن یک فاز ترکیب داده‌ها وجود داشته باشد. در این روش‌ها زمان ترکیب و ارتباط عامل‌ها مشخص و ثابت است و در روش‌هایی دیگر چون پند دهی عامل‌ها هر زمان که نیاز به کمک داشته باشند می‌توانند از عامل‌های دیگر بهره ببرند. باید تأکید کرد که روش‌های ارائه‌شده در این پژوهش همانند روش WSS زمان ترکیب داده‌های ثابتی در نظر گرفته شده است.

## فصل سوم

### مفاهیم علمی پیش نیاز پایان نامه

#### ۳-۱ مقدمه

در این فصل سعی شده تا موضوعاتی که در روش پیشنهادی به کار رفته اند و به درک بهتر موضوع کمک می کنند تشریح می شوند در این جهت ابتدا در مورد روش های یادگیری و روش یادگیری  $Q$  که معمولاً در کارهای یادگیری مشارکتی استفاده میشود توضیح داده شده است. مطمئناً شناخت یادگیری تقویتی حتی به صورت جزئی می تواند در درک یادگیری مشارکتی بسیار مؤثر باشد. از آنجایی که بررسی عملکرد روش پیشنهادی با روش یادگیری مشارکتی بر مبنای کوتاه ترین فاصله تجربه شده انجام می شود در ادامه به تشریح معیارهای SEP و شک پرداخته خواهد شد. بعد از آن محیط های آزمایشی و معیارهای ارزیابی استفاده شده در آزمایش های این پژوهش تشریح خواهد شد. همچنین از آنجایی که در این پژوهش از انتگرال فازی چوکت استفاده شده است گذری خلاصه بر اندازه گیری های فازی و غیرافزایشی شده است و سپس از بین انتگرال های فازی دو انتگرال همه کاره سوگنو و چوکت که می توان به روی هر نوع داده ای اعمال کرد را معرفی کردیم و نشان داده شده است چرا در کاربرد مورد استفاده در این پژوهش فقط از انتگرال چوکت بهره برده شده است.

---

```

1: procedure Q-LEARNING
Ensure: Initialize the  $Q$  matrix;
2:   while not End Of Learning do
3:     Visit the state  $s$ ;
4:     Select an action  $a$  based on an action selection policy;
5:     Carry out the  $a$  and observe a reward  $r$  at the new state  $s'$ ;
6:      $Q[s, a] \leftarrow Q[s, a] + \alpha(r + \lambda \max_{a'}(Q[s', a']) - Q[s, a])$ ;
7:      $s \leftarrow s'$ ;
8:   end while
9:   return  $Q$ ;
10: end procedure

```

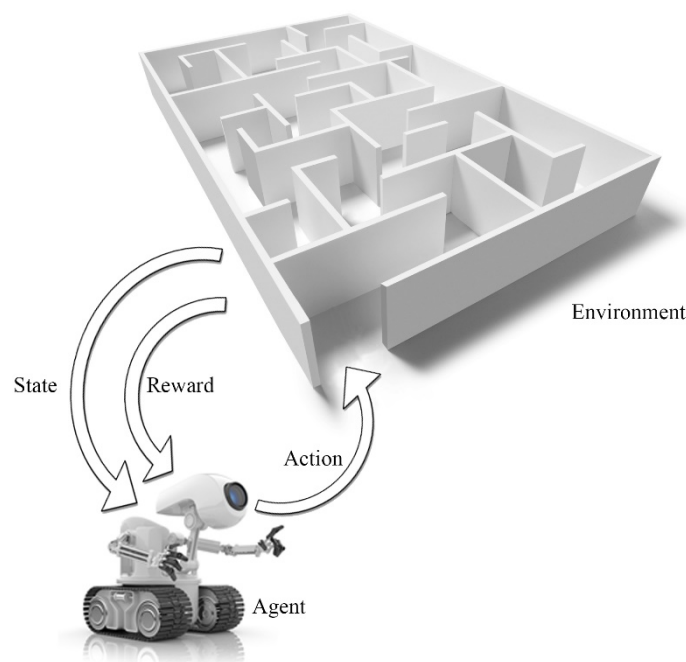
---

### ۳-۲ یادگیری تقویتی

یادگیری تقویتی معمولاً به آن سری از روش‌های یادگیری گفته می‌شود که عامل به دنبال رسیدن به یادگیری از طرق ارتباط با محیط است. در این دسته از روش‌ها به عامل ارزش اعمال گفته نمی‌شود و عامل با تعاملی که با محیط دارد باید ارزش اعمال را کشف نماید. در این روش فرایند یادگیری به بخش‌هایی با عنوان چرخه یادگیری شکسته می‌شود، که الگوریتم ۳-۱ به ازای هر چرخه یادگیری به اجرا در می‌آید. هر چرخه یادگیری از قرار دادن عامل در یک حالت تصادفی شروع و تا رسیدن به یک حالت پایانی ادامه دارد. در طول هر چرخه تا زمان رسیدن به حالت پایانی عامل وظیفه دارد عمل را انتخاب نماید و پس از دریافت پاداش عمل انجام‌شده به بروز رسانی اطلاعات پردازد. بروز رسانی داده‌ها بر اساس فرایند تصادفی مارکف و روش برنامه نویسی پویا می‌تواند انجام شود. بر اساس تئوری ارائه‌شده در فرایند تصادفی مارکف باید انتخاب عمل به صورتی باشد که انجام هر عمل در هر حالت ضمانت شود. معمولاً در یادگیری تقویتی این انتخاب عمل به وسیله روش‌هایی چون بولتزمن انجام می‌شود. در شکل ۳-۱ می‌توان فرایند یادگیری را مشاهده کرد.

### ۳-۳ روش‌های انتخاب عمل

همان‌طور که گفته شد یکی از وظایف عاملی که از یادگیری تقویتی استفاده می‌کند انتخاب عمل در طول یادگیری است. عامل بعد از رسیدن به هر حالت تا زمان رسیدن به یک حالت پایانی باید اعمالی انتخاب نماید. این انتخاب می‌تواند کاملاً بر اساس داده‌های جمع‌آوری شده در چرخه‌های یادگیری قبلی باشد که اصطلاحاً بهره‌برداری نامیده می‌شود. برای این کار کافی است در هر حرکت عملی انتخاب شود که در جدول  $Q$  ارزش بالاتری دارد اما این کار باعث می‌شود تا عامل مسیرهایی را تکرار نماید و شرط بررسی تمام مسیرها که از فرایند تصادفی مارکو به یادگیری تقویتی رسیده است را ارضا نمی‌کند. پس عامل نیاز است در طول یادگیری گاهی



شکل ۳-۱: شمایی از فرایند یادگیری تقویتی در تعامل با محیط [۲]

فارغ از بهترین عمل اعمال دیگر را نیز بررسی نماید. که این بررسی اعمال دیگر را اکتشاف گویند. در بدترین حالت می‌توان گفت که در طول یادگیری عامل به صورت تصادفی اعمال را انتخاب نماید اما این موضوع فرایند یادگیری را بسیار طولانی می‌کند. پس عامل باید برای رسیدن به یک تعادل در اکتشاف و بهره‌برداری یک روش مناسب در انتخاب اعمال داشته باشد. در ادامه دو روش انتخاب عمل که در یادگیری مشارکتی مورد استفاده قرار می‌گیرد خواهد آمد.

### ۳-۳-۱ $\epsilon$ -حریصانه

در  $\epsilon$ -حریصانه جهت رسیدن به یک تعادل در اکتشاف و بهره‌برداری یک پارامتر  $\epsilon$  که مقداری بین ۰ و ۱ است در نظر گرفته میشود. سپس به احتمال  $\epsilon$  اکتشاف و به احتمال  $(1 - \epsilon)$  بهره‌برداری از داده‌های جدول  $Q$  انجام می‌شود. جهت پیاده‌سازی این روش یک مقدار تصادفی  $\xi$  بین ۰ و ۱ تولید شده و در صورتی که این  $\xi$  کوچک‌تر از  $\epsilon$  باشد حرکت تصادفی و در غیر این صورت بهترین حرکت بر اساس داده‌های جدول  $Q$  انجام میشود.

### ۳-۳-۲ بولتزمن

روش  $\epsilon$ -حریصانه توانسته تا حدودی بین اکتشاف و بهره‌برداری تعادل ایجاد نماید اما احتمال انتخاب اعمال در  $\epsilon$ -حریصانه ثابت بوده و ارتباطی با ارزش اعمال در جداول  $Q$  ندارد. روش بولتزمن با بهره‌گیری از رابطه ۳-۱ سعی دارد تا احتمال انتخاب عمل  $i$ ام در یک موقعیت  $s$  را بر اساس ارزش اعمال در همان موقعیت در جدول  $Q$  محاسبه شود [۲۰]. همچنین در این رابطه مشخص است این روش دارای پارامتری به‌عنوان  $\tau$  برای



جدول ۳-۱: ساختار جدول CP [۲]

...	آخرین عمل انجام شده در حالت $i$	...
...	حالت بعدی که مشاهده شده	...
...	شماره آخرین گام حرکت مشاهده حالت $i$	...

کنترل حساست به اختلاف ارزش بین اعمال در نظر گرفته شده است هرچقدر این میزان بزرگتر باشد به اختلاف ارزشها اهمیت کمتری می شود. یعنی در صورتی که  $\tau \rightarrow \infty$  تمامی اعمال ممکن در موقعیت  $s$  به احتمال یکسانی انتخاب می شود و در صورتی که  $\tau \rightarrow 0$  به صورت حریصانه عملی انتخاب می شود که مقدار  $Q$  بیشتری دارد.

$$p_Q(a|s) = \frac{e^{\frac{Q[s,i]}{\tau}}}{\sum_b e^{\frac{Q[s,b]}{\tau}}} \quad (۱-۳)$$

### ۳-۴ الگوریتم مورد مقایسه با روش پیشنهادی

در فصل دوم به تشریح روش یادگیری مشارکتی بر مبنای «کوتاه ترین فاصله تجربه شده» پرداخته شد اما از آنجایی که بررسی عملکرد روش پیشنهادی در این پژوهش با کار میرزایی [۲] انجام شده است که مدرن ترین روش ارائه شده در این زمینه می باشد، لذا لازم است شناخت بیشتری از معیارهای به کاررفته در این کار صورت گیرد.

#### ۳-۴-۱ معیار کوتاه ترین فاصله تجربه شده

میرزایی معیاری با عنوان «کوتاه ترین فاصله تجربه شده» تعریف کرد [۲]، اما بهتر بود «کوتاه ترین فاصله استدلال شده» نامیده شود چراکه در روش محاسبه ارائه توسط ایشان کوتاه ترین مسیرهایی را که بتوان از مسیرهای تجربه شده استدلال کرد پیدا می نماید.

این معیار برای هر عمل از هر حالت مقداری در نظر می گیرد که در پایان هر چرخه یادگیری بروز رسانی می شود. جهت به روز رسانی این جدول یک ماتریس با نام  $CP^1$  در نظر گرفته شده است که وظیفه نگهداری مسیری طی شده در هر چرخه را بر عهده دارد. این ماتریس به تعداد حالت های محیط ستون دارد و سه سطر دارد. مطابق با جدول ۳-۱ در سلول  $(1, i)$  آخرین عمل انجام شده در حالت  $i$  انجام می شود، در سلول  $(2, i)$  حالت مقصدی که عامل بعد از این حالت مشاهده کرده و در سلول  $(3, i)$  شماره آخرین گام حرکت در چرخه یادگیری فعلی که عامل در حالت  $i$  بوده نگهداری می شود.

<sup>1</sup>Current Path

بعد از اتمام هر چرخه یادگیری زمان بروز رسانی جدول *SEP* با بهره‌گیری از جدول *CP* است. در الگوریتم ارائه شده توسط میرزایی [۲] ابتدای سلول مربوط به عملی که در آخرین چرخه یادگیری انجام شده مقدار ۱ می‌گیرد چراکه انجام این حرکت باعث رسیدن به نقطه هدف شده است پس فاصله این حرکت تا مقصد برابر با یک خواهد بود. بعد از آن به ترتیب معکوس گام‌های یادگیری سلول مربوط به اعمال انجام شده بروز رسانی می‌شود. جهت بروز رسانی مقدار یک خانه برابر حداقل مقدار مربوط به اعمال مقصد در همین جدول *SEP* به علاوه ۱ خواهد بود. چرا که فاصله حالت جاری تا حالت بعدی ۱ بوده و جمع این فاصله با کوتاه‌ترین فاصله در حالت مجاور کوتاه‌ترین فاصله حالت جاری را خواهد ساخت.

### ۳-۴-۲ شوک

بزرگ‌ترین دلیل رسیدن به یادگیری در یادگیری تقویتی پاداش‌های دریافتی برای انجام اعمال از محیط است. اما هرچه این اعمال به حالت پایانی نزدیک‌تر باشند باز ارزش‌تر هستند. میرزایی در [۲] معتقد است که پاداش دریافتی برای رسیدن به حالت‌های نهایی دارای ارزش بالاتری است و عامل حالتی که اثر بیشتری از این پاداش گرفته باشد بیشتر می‌شناسد. بر همین اساس در پارامتر شوک تعداد دفعاتی که اثر پاداش دریافتی از حالت‌های نهایی به هر خانه از جدول *Q* رسیده است شمارش می‌شود. برای شمارش این کار کافی است در شروع یادگیری یک ماتریس در ابعاد ماتریس *Q* با مقادیر اولیه ۰ تولید کرده و پس از انجام هر عمل در صورتی که حالت بعدی حالت پایانی باشد یا حالتی باشد که قبلاً از پاداش حالت پایانی اثر گرفته باشد مقدار مربوط به آن عمل در جدول شوک یک واحد افزایش می‌یابد.

### ۳-۵ محیط‌های آزمایش

راسل در [۲۱] محیط‌ها را با پنج دیدگاه دسته‌بندی می‌نماید.

- **مشاهده پذیر و نیمه مشاهده پذیر:** اگر عامل به کمک حسگرهای خود توانایی تشخیص حالت محیط را داشته باشد محیط را مشاهده پذیر گویند. در غیر این صورت محیط نیمه مشاهده پذیر یا تا حدی قابل مشاهده نامیده می‌شود.
- **قطعی و غیرقطعی:** اگر بتوان حالت بعدی محیط را بر اساس سابقه اعمال و حالت فعلی مشخص کرد محیط قطعی است و در غیر این صورت محیط غیرقطعی است.
- **واقع‌ای و غیرواقع‌ای:** در صورتی که هر مرحله از مراحل دیگر مستقل باشد محیط را دوره‌ای می‌نامیم.
- **ایستا و پویا:** اگر محیط در مدت زمان بین درک و انتخاب عمل تغییر کند پویا و در غیر این صورت ایستا

					G
G					
					G

شکل ۳-۲: محیط پلکان مارپیچ [۲]

است.

• **گسسته و پیوسته:** اگر مشاهدات و اعمال به شکل جداگانه تعریف شوند محیط را گسسته گویند.

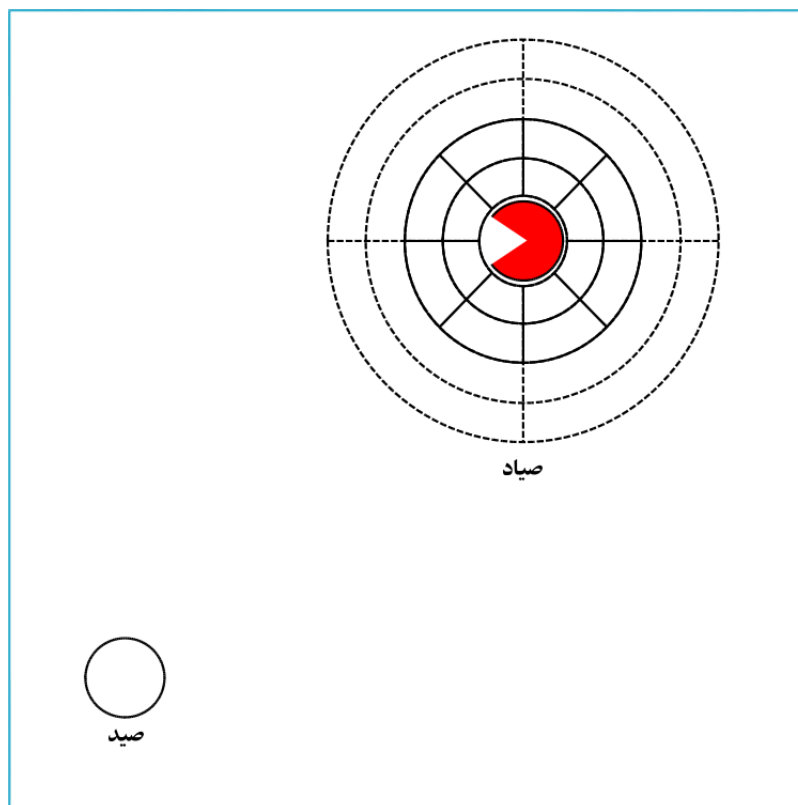
معمولا بعد از ارائه‌ی هر روش یادگیری نیاز است تا عملکرد آن روش در محیط‌های مختلف مورد ارزیابی قرار گیرد. در این پژوهش نیز دو محیط جهت این امر مورد استفاده قرار گرفته است که در ادامه تشریح میشوند.

### ۳-۵-۱ محیط پلکان مارپیچ

پلکان مارپیچ [۱، ۲] همان‌طور که در شکل ۳-۲ مشخص است یک محیط ایستا است که از یک مربع  $6 \times 6$  شامل ۳ خانه هدف، تعدادی دیوار و تعدادی خانه آزاد تشکیل شده است. عامل در این محیط باید سیاست رسیدن از هر حالت به یک حالت هدف را با استفاده از چهار عمل اصلی بالا، پایین، چپ، راست را یادگیری نماید. در صورتی که انجام عمل توسط عامل باعث انتقال عامل به خانه هدف شود به عامل پاداش ۱۰، در صورتی که عمل انتخابی عامل را به دیوار بزند پاداش ۱-، در غیر این صورت معکوس فاصله حالت جاری تا نزدیک‌ترین هدف را به عنوان پاداش دریافت می‌نماید.

### ۳-۵-۲ محیط صید و صیاد

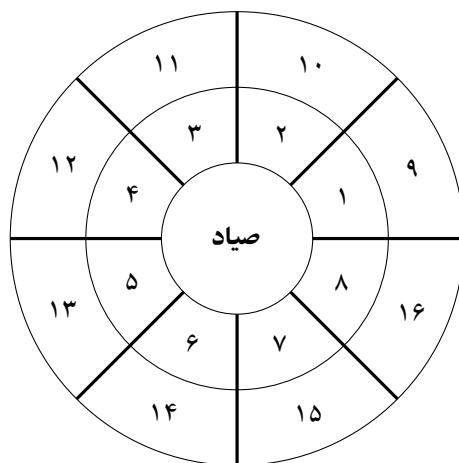
همان‌طور که در شکل ۳-۳ محیط صید و صیاد آورده شده است، محیطی است شامل یک مربع  $10 \times 10$  شامل دو عامل صید و صیاد که در پیاده‌سازی‌ها معمولا عامل صید به صورت تصادفی حرکت کرده عامل صیاد روش شکار را یادگیری می‌نماید [۱، ۲]. این محیط برخلاف محیط پلکان مارپیچ یک محیط پیوسته است که عامل‌ها می‌توانند در هر نقطه از آن قرار گیرند. حرکت عامل‌ها در این محیط هم به صورت پیوسته است به این صورت که عامل صیاد می‌تواند به هر نقطه با شعاع یک و عامل صید به هر نقطه با شعاع ۰,۵ در اطرافشان حرکت



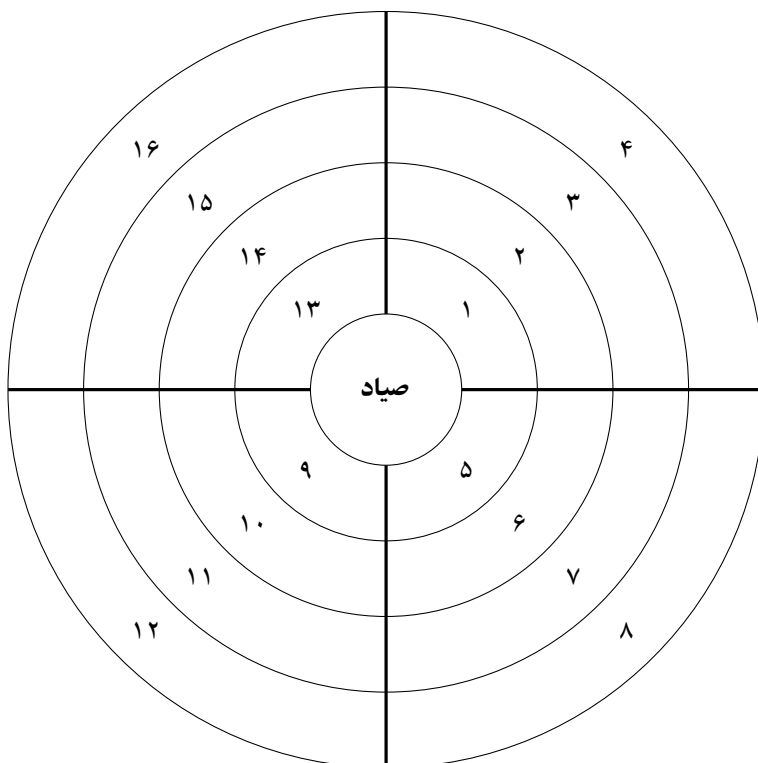
شکل ۳-۳: محیط صید و صیاد

نماید. اما از آنجایی که روش‌های یادگیری تقویتی نیاز به تعداد اعمال مشخص دارد باید این پیوستگی گسسته سازی شود. در اینجا برای گسسته سازی مطابق شکل ۳-۴ (آ) زاویه حرکت به فاصله‌های ۴۵ درجه‌ای و فاصله حرکت به نیم و یک تقسیم شده است که جمعاً تولید ۱۶ عمل برای عامل می‌نماید.

موضوع بعد حالت‌های قرارگیری عامل است. در محیط صید و صیاد برای صیاد یک دامنه دید در نظر گرفته می‌شود که در صورتی که عامل صید در فاصله کمتر مساوی از دامنه دید صیاد باشد صیاد قادر به دیدن صید خواهد بود. از آنجایی که این محل قرارگیری عامل صید نسبت به صیاد یک مقدار پیوسته است نیاز به گسسته سازی دارد. در کار پیش رو دامنه دید عامل صیاد برابر ۲ در نظر گرفته شده است و برای گسسته سازی این دامنه دید را مطابق شکل ۳-۴ (ب) به زوایای ۹۰ درجه و فاصله‌ی به اندازه‌های ۰/۵ تقسیم می‌نماییم. پس تعداد حالت‌ها زمانی که صید در دامنه دید قرار گرفته باشد برابر با ۱۶ حالت خواهد بود این شانزده به علاوه یک حالت که عامل در دامنه دید نباشد ۱۷ حالت را برای سیستم به وجود می‌آورد. همان طور که گفته شد حرکت‌های صید نیز در این سیستم به صورت تصادفی انجام می‌شود. عامل صیاد در صورتی که با انجام یک عمل صید را شکار کند پاداش ۱۰ و در غیر این صورت پاداش ۰/۱ - دریافت می‌نماید.

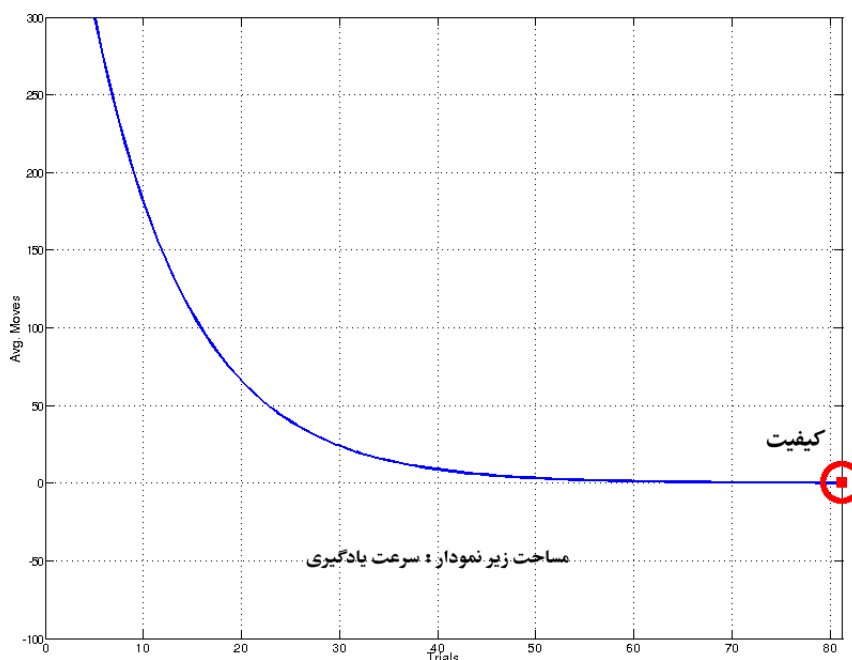


(آ) ۱۶ عمل تعریف شده برای صیاد



(ب) دامنه دید تعریف شده برای صیاد

شکل ۳-۴: دامنه‌ی دید و حالت تعریف شده برای عامل صیاد در محیط صید و صیاد [۲].



شکل ۳-۵: سرعت و کیفیت یادگیری از معیارهای ارزیابی و مقایسه‌ی عملکرد الگوریتم‌های یادگیری تقویتی می‌باشد [۲].

### ۳-۶ معیارهای ارزیابی

ارزیابی هر سیستمی نیاز به معیارهایی دارد، که در یادگیری مشارکتی نیز دو معیار «سرعت و دقت یادگیری» وجود دارد که بر اساس میانگین تعداد قدم‌ها در هر چرخه یادگیری محاسبه می‌شوند [۱، ۲]. در صورتی که تعداد قدم‌های چرخه‌های یادگیری نگهداری و به صورتی نموداری که در محور عمودی میانگین تعداد قدم‌ها و در محور افقی چرخه یادگیری باشد رسم شود، می‌توان آخرین نقطه از محور را که نشانه میانگین تعداد چرخه‌ها در آخرین چرخه یادگیری است را به عنوان کیفیت و مساحت زیر این نمودار را به عنوان سرعت یادگیری در نظر گرفت که هدف کمینه کردن این معیارها است؛ در شکل ۳-۵ این معیارها نمایش داده شده است.

### ۳-۷ اندازه‌گیری و انتگرال فازی

برای درک روش پیشنهادی نیاز به داشتن اطلاعات پایه در مورد اندازه‌گیری‌های فازی<sup>۱</sup> و انتگرال فازی که با هدف جمع‌آوری اطلاعات<sup>۲</sup> ارائه شده‌اند، داریم. اندازه‌گیری‌های فازی پیش‌زمینه‌ای بر انتگرال‌های فازی هستند که قبل از آشنایی با انتگرال‌های فازی نیاز به معرفی اندازه‌گیری‌های فازی داریم. اگر فرض کنیم که تعدادی منبع اطلاعاتی  $X = \{x_1, x_2, \dots, x_n\}$  که این منابع اطلاعاتی دریافتی از حسگرها، پاسخ‌های داده شده به یک پرسشنامه و غیره وجود داشته باشند. اندازه‌گیری فازی میزان ارزش اطلاعاتی این منابع را در اختیار ما می‌گذارد. معمولاً اندازه‌گیری فازی توسط تابع  $g : 2^{|X|} \rightarrow [0, 1]$  تعریف می‌شود که ورودی آن یک

<sup>۱</sup>Fuzzy measures

<sup>۲</sup>Aggregate Information

زیرمجموعه‌ای از منابع اطلاعاتی می‌باشد و خروجی آن یک مقدار مابین صفر و یک که میزان ارزش اطلاعاتی که آن زیرمجموعه از منابع اطلاعاتی ورودی تابع را مشخص می‌کند. این تابع باید دارای شرایط مرزی تعریف شده و یکنواختی باشد که در ادامه به معرفی شرایط می‌پردازیم [۲۲]:

۱. شرایط مرزی: اگر اطلاعاتی در دست نداریم ارزش صفر را دارد و کلیه اطلاعات ارزش ۱ را دارد.

$$g(\emptyset) = 0, \quad g(X) = 1 \quad (۲-۳)$$

۲. یکنواختی - غیر کاهشی: اگر اطلاعات بیشتری به دست آمد ارزش کلیه اطلاعات که شامل اطلاعات جدید می‌باشد حداقل به اندازه زمانی است که آن اطلاعات جدید بدست نیامده است.

$$A \subseteq B \subseteq X \Rightarrow g(A) \leq g(B) \leq 1 \quad (۳-۳)$$

مقادیر تابع  $g$  یا توسط کارشناس ارائه می‌شود یا توسط یک تابعی مدل می‌شود، یکی از توابع معروف برای تخمین مقادیر تابع  $g$  تابع اندازه‌گیری- $\lambda$  سوگنو<sup>۱</sup> می‌باشد که به صورت زیر تعریف می‌شود [۲۳].

$$g(\{x_1, \dots, x_l\}) = \frac{1}{\lambda} \left[ \prod_{i=1}^l (1 + \lambda g_i) - 1 \right] \quad (۴-۳)$$

که در معادله ۴-۳ مقدار  $g_i$  ها مقادیر ارزش هریک از منابع اطلاعاتی است و  $\lambda$  بگونه‌ای تعیین می‌گردد که  $g(X) = 1$  شود که این مقدار برابر با جواب معادله‌ی زیر باشد.

$$\lambda + 1 = \prod_{i=1}^n (1 + \lambda g_i), \quad \lambda \in (-1, \infty) \quad (۵-۳)$$

نکته‌ای که در رابطه با تابع اندازه‌گیری- $\lambda$  سوگنو باید توجه کرد این است که به ازای مقادیر  $n$  مختلف باید ریشه‌یابی بروی متغیر  $\lambda$  صورت گیرد؛ این ویژگی باعث می‌شود که این تابع در بعضی از کاربردها کارایی نداشته باشد. به عنوان مثال در کار این پژوهش از آنجایی که تعداد عامل‌ها متفاوت می‌باشد، در صورت استفاده از تابع اندازه‌گیری- $\lambda$  باید به ازای هر دفعه تغییر در تعداد عامل‌ها (از آنجایی که مقادیر  $n$  تغییر می‌کند) یکبار روی متغیر  $\lambda$  ریشه‌یابی صورت بگیرد که امکان همچنین ریشه‌یابی‌ای بدون سربار محاسباتی سنگین میسر نیست. لذا در کاربردهایی که تعداد  $n$  متغیر می‌باشد استفاده از تابع اندازه‌گیری- $\lambda$  عاقلانه نیست.

<sup>۱</sup> Sugeno  $\lambda$ -Measure

انتگرال فازی در واقع یک تعمیمی به روش میانگین وزنی<sup>۱</sup> می باشد بطوری که نه تنها مشخصه های مهم تک تک ویژگی ها را در نظر می گیرد بلکه اطلاعات تعاملات بین ویژگی ها را نیز در نظر می گیرد [۲۴]. از میان انتگرال های فازی دو انتگرال سوگنو<sup>۲</sup> و چوکت<sup>۳</sup> از الگوریتم هایی هستند که می توانند بر روی هر اندازه گیری فازی مورد استفاده واقع شوند [۲۵]. فرض کنیم که تابعی چون  $h : X \rightarrow [0, 1]$  وجود دارد که مقادیر منابع اطلاعاتی را به بازه ی  $[0, 1]$  نگاشت می کند. در واقع  $h$  تابع پشتیبان<sup>۴</sup> منابع اطلاعاتی می باشد. انتگرال فازی سوگنو به صورت ۶-۳ تا ۸-۳ تعریف می شود [۲۵، ۲۶]:

$$\int_s h \circ g = \mathcal{S}_g(h) = \bigvee_{i=1}^n h(x_{\pi_i^s}) \wedge g(A_i^s) \quad (۶-۳)$$

$$h \xrightarrow{\pi^s} h(x_{\pi_1^s}) \leq h(x_{\pi_2^s}) \leq \dots \leq h(x_{\pi_n^s}) \quad (۷-۳)$$

$$A_i^s = \{x_{\pi_i^s}^s, x_{\pi_{i+1}^s}^s, \dots, x_{\pi_n^s}^s\} \quad (۸-۳)$$

در روابط ۶-۳ تا ۸-۳ نماد  $s$  بیان گر «سوگنو» می باشد. در انتگرال سوگنو لازم است که مقادیر منابع اطلاعاتی را مرتب کنیم که  $\pi^s$  عملگر جایگشت انتگرال فازی سوگنو می باشد که خروجی مقادیر تابع  $h$  را به ترتیب صعودی مرتب می کند. نمادهای  $\vee$  و  $\wedge$  به ترتیب عملگرهای  $\max$  و  $\min$  می باشد. در این انتگرال ابتدا مقادیر دریافتی از منابع اطلاعاتی به تابع پشتیبان  $h$  ارسال می شود و سپس مقادیر خروجی تابع پشتیبان به ازای همه ی اطلاعات دریافتی را به صورت صعودی توسط عملگر جایگشت  $\pi^s$  مرتب می شود. مجموعه ی  $A_i^s$  اندیس عناصر مرتب شده مقادیر تابع پشتیبان از اندیس  $i$ ام تا اندیس  $n$  می باشد. سپس طبق آنچه که در ۶-۳ آمده است از کوچکترین مقدار  $h(x_{\pi_1^s})$  شروع می کنیم با  $g(A_1^s)$  کمینه گیری می کنیم و سپس می رویم به دومین کوچکترین عنصر  $h(x_{\pi_2^s})$  و همین کار را تا آخرین (بزرگترین) عنصر تکرار می کنیم و سپس یک بیشینه گیری روی این مقادیر انجام می دهیم که خروجی انتگرال سوگنو می شود.

انتگرال فازی چوکت به صورت ۹-۳ تعریف می شود [۲۵، ۲۷]. در این رابطه  $f : X \rightarrow \mathbb{R}$  می باشد که از وجه تمایز انتگرال فازی چوکت با سوگنو می باشد و  $\pi^c$  عملگر جایگشت انتگرال فازی چوکت می باشد.

<sup>1</sup>Weighted Arithmetic Mean

<sup>2</sup>Sugeno

<sup>3</sup>Choquet

<sup>4</sup>Support



$$\int_c f \circ g = C_g(f) = \sum_{i=1}^n \left( f(x_{\pi_i^c}) - f(x_{\pi_{i-1}^c}) \right) \cdot g(A_i^c) \quad (9-3)$$

$$f \xrightarrow{\pi^c} f(x_{\pi_1^c}) \leq f(x_{\pi_2^c}) \leq \dots \leq f(x_{\pi_n^c}) \quad (10-3)$$

$$A_i^c = \{x_{\pi_i^c}^c, x_{\pi_{i+1}^c}^c, \dots, x_{\pi_n^c}^c\} \quad (11-3)$$

$$f(x_{\pi_0^c}) = 0 \quad (12-3)$$

در روابط ۹-۳ تا ۱۱-۳ نماد  $c$  بیانگر «چوکت» می‌باشد. عملکرد انتگرال چوکت شباهت نزدیکی با انتگرال سوگنو دارد به این صورت که در انتگرال چوکت مقادیر دریافتی از منابع اطلاعاتی را به تابع پشتیبان  $f(\cdot)$  ارسال می‌شود و خروجی این تابع را به ازای تمامی ورودی‌ها توسط عملگر جایگشت  $\pi^c$  به صورت صعودی مرتب می‌کنیم. تعریف  $A_i^c$  در ۱۱-۳ مشابه انتگرال سوگنو می‌باشد. طبق تعریف انجام شده در ۱۲-۳ مقدار ۰ را به اول مقادیر مرتب شده  $f(\cdot)$  در ۱۰-۳ اضافه می‌کنیم. سپس توسط رابطه‌ی ۹-۳ مجموع ضرب اختلاف دو عنصر متوالی مرتب شده  $f(\cdot)$  در  $A_i^c$  را به عنوان خروجی انتگرال چوکت حساب می‌کنیم.

انتگرال‌های فازی سوگنو و چوکت در حالت کلی دارای تفاوت‌هایی هستند که از جمله‌ی مهم‌ترین این ویژگی‌ها تفاوت تعریف توابع  $h$  و  $f$  در این انتگرال‌ها می‌باشد که باعث می‌شود انتگرال چوکت برای مسائلی که مقادیر اعداد حائز اهمیت است، مناسب باشد و از طرف دیگر انتگرال سوگنو زمانی مطلوب است که تنها ترتیب اعداد مد نظر باشد [۲۸]. به همین علت در این پژوهش انتگرال فازی چوکت مورد استفاده قرار گرفته است زیرا که ورودی انتگرال اعداد کاملاً معنی‌دار می‌باشد و اعمال تابع  $h$  بروی مقادیر منابع اطلاعاتی، معانی آن‌ها را تغییر داده و اطلاعات نامطلوبی تولید خواهد کرد.

### ۸-۳ نتیجه‌گیری

در این فصل مطالبی که برای درک بهتر روش پیشنهادی نیاز است شرح داده شده است. در روش پیشنهادی از یادگیری  $Q$  با سیاست‌های انتخاب عمل بولترمن،  $\varepsilon$ -حریصانه و انتگرال فازی چوکت استفاده شده است. در ادامه، در فصل بعدی روش پیشنهادی توضیح داده شده است.

## فصل چهارم

### روش پیشنهادی

#### ۴-۱ مقدمه

در این فصل جزییات روش پیشنهادی به طور مفصل معرفی خواهد شد، روش ارائه شده در حالت کلی از دو قسمت تشکیل شده است؛ اولین و مهم‌ترین قسمت ارائه یک معیار خبرگی جدید به نام معیار خبرگی «ارجاع» که برای هر عامل در هر چرخه یادگیری محاسبه و در یک «ماتریس ارجاع» نگه‌داری می‌شود. دومین قسمت مربوط به ترکیب دانش‌های عامل‌ها هستند که با استفاده از یک مدل انتگرال فازی، صورت می‌گیرد. همانطور که در فصل بعدی نیز نشان داده خواهد شد استفاده از مدل انتگرال فازی به دلیل خواصی مهمی که این مدل دارد باعث می‌شود سرعت و کیفیت یادگیری به طرز چشم‌گیری افزایش یابد. در این فصل ابتدا به معرفی معیار «ارجاع» و دلیل استفاده از آن می‌پردازیم سپس یادگیری مشارکتی چندعامله با استفاده از ماتریس ارجاع و انتگرال فازی معرفی خواهد شد و در نهایت نشان داده خواهد شد که استفاده از انتگرال فازی می‌تواند نتایج بهتری را نسبت به مدل‌های سنتی چون مدل مجموع وزنی<sup>۱</sup> ارائه دهد.

---

<sup>۱</sup> Weighted Sum

## ۲-۴ معیار خبرگی - ماتریس ارجاع و خاطره

در دنیای واقعی «خبرگی» تعاریف متعددی به خود گرفته است، در روانشناسی خبرگی به معنی عملکرد برتر عامل تلقی می‌شود. در جامعه شناسی خبره به فردی گفته می‌شود که برچسب خبرگی توسط یک گروهی به فرد زده شده است و آن گروه به توانایی که آن فرد در اختیار دارد علاقه‌مند<sup>۱</sup> است. در فلسفه خبره به فردی گفته می‌شود که دانشی که فرد تازه‌کار در اختیار ندارد را دارا می‌باشد [۲۹]. اگر تعاریف مختلف «خبرگی» را بررسی کنیم می‌بینیم که همه‌ی تعاریف در واقع تعبیری از میزان کیفیت عملکرد عامل نسبت به دیگر عامل‌ها می‌باشد. این تعبیر کلی از «خبرگی» انگیزه‌ای شد که درصدد معرفی معیاری برآییم که در حالت کلی بتوان به کلیه‌ی تعاریف «خبرگی» قابل تعمیم باشد.

**فرضیه ۱-۴ (خبرگی).** فرض می‌کنیم عامل  $A$  در محیط  $\mathcal{E}$  در پی رسیدن به یک مجموعه اهداف  $G \subseteq \{g_1, g_2, \dots, g_n\}$  می‌باشد. میزان خبرگی عامل رابطه‌ی معکوسی با میزان تلاش عامل برای رسیدن به اهداف تعریف شده خود دارد.

طبق آنچه که در فرضیه بالا آورده شده است از بین چند عاملی که در یک محیط و یک مجموعه از اهداف فعالیت می‌کنند، عاملی خبره‌تر است که تلاش کمتری برای رسیدن به آن مجموعه اهداف می‌کند. شاید این مساله در نگاه اول نامتعارف به ذهن برسد ولی در فعالیتهای روزمره ما انسان‌ها نیز به کرات شاهد این امر می‌باشیم. به عنوان مثال رانندگی دو فرد مبتدی و حرفه‌ای را در نظر بگیریم؛ فرد مبتدی هنگام رانندگی تمام حواس خود را معطوف به رانندگی می‌کند تلاش بسیار زیادی برای کنترل نسبت میزان کلاچ و گاز می‌کند و هنگام رانندگی به طور طبیعی رانندگی نمی‌کند و ... ولی فرد خبره کلیه موارد ذکر شده را بطور خودکار و طبیعی انجام می‌دهد بطوری که انگار رانندگی مانند دیگر رفتارهای طبیعی وی چون نفس کشیدن می‌باشد، که بصورت خودکار صورت می‌پذیرد. از این گونه مثال‌ها از کاربرد فرضیه ۱-۴ در زندگی روزمره ما زیاد می‌توان یافت.

توجه شود که در فرضیه ۱-۴ عبارت «میزان تلاش» عامل می‌تواند در کاربردهای مختلف تعبیر مختلفی به خود بگیرد، مثلاً در مثال راننده‌ی مبتدی و خبره میزان نسبت مسافت طی شده بر زمان رانندگی را می‌توان به عنوان «میزان تلاش» عامل در نظر گرفت که در شرایط یکسان راننده‌ی خبره‌تر به طور نسبی در زمان کوتاه‌تری یک مسافت مشخصی را طی خواهد کرد (در رد کردن پیچ و خم‌های ترافیک و مدت زمان ترمز و ... زمان کمتری را تلف می‌کند). یا به عنوان مثال دیگر، دانشجوی قوی و دانشجوی ضعیف را مورد بررسی قرار دهیم، دانشجویی خبره هست که زمان کمتری را صرف حل صحیح یک مساله خاص کند (با فرض اینکه دانشجویها حتماً باید مساله را حل کنند). همانطور که دیدیم کمیت «میزان تلاش» عامل برای مسائل مختلف معیار متفاوتی را دربر می‌گیرد ولی همگی از همان اصل معرفی شده در فرضیه ۱-۴ تبعیت می‌کنند.

در یادگیری مشارکتی با استفاده از فرضیه ۱-۴ می‌توان با تعریف ۱-۴ یک معیار خبرگی جدید را معرفی

کرد که مبنا و پایه‌ی دستاوردهای این پژوهش می‌باشد.

**تعریف ۴-۱** (معیار خبرگی «میزان ارجاع»). فرض می‌کنیم مجموعه‌ای از عامل‌ها  $\mathbb{A} = \{A_1, A_2, \dots, A_m\}$  در محیط  $\mathcal{E}$  در پی رسیدن به یک مجموعه اهداف  $\mathcal{G} \subseteq \{g_1, g_2, \dots, g_n\}$  می‌باشند. اگر ما به طور مجازی و دلخواه محیط  $\mathcal{E}$  را به  $k$  ناحیه مانند  $e_i$  افراز کنیم بطوری که  $\mathcal{E} = \{e_i \mid \bigcup_{i=1}^k e_i = \mathcal{E} \wedge \forall i, j \in \{1, 2, \dots, k\} \wedge i \neq j : e_i \cap e_j = \emptyset\}$  میزان ارجاع هر عامل در هر ناحیه را میزان حضور آن عامل را در آن ناحیه تعریف می‌کنیم.

در تشریح آنچه که در تعریف ۴-۱ آمده است می‌توان گفت که در سیستم‌های چندعاملی که همگی عوامل در یک محیط به صورت مستقل در حال فعالیت هستند؛ محیط را به چند ناحیه دلخواه افراز می‌کنیم که اجتماع نواحی باهم کل محیط  $\mathcal{E}$  را تشکیل دهند و هیچ دو ناحیه‌ای اشتراکی باهم نداشته باشند [۳۰]. در این چنین افرازی از محیط، در هر ناحیه عاملی که نسبت به بقیه خبره‌تر است، نسبت به بقیه عوامل در همان ناحیه میزان تمایل حضور کمتری را از خود نشان می‌دهند. به عبارت دیگر عاملی که خبره‌تر است تمایل دارد کوتاه‌ترین مسیر رسیدن به اهداف خود را طی کند که نهایتاً منجر خواهد شد که میزان حضور عامل در هریک از نواحی محیط کمینه شود.

آنچه که در فرضیه ۴-۱ در مورد «میزان تلاش» عامل آمده است در تعریف ۴-۱ در به صورت «میزان حضور عامل در هر ناحیه» تعریف شده است. بطوری که طبق فرضیه مطرح شده میزان خبرگی عامل در هر ناحیه رابطه‌ی معکوسی با میزان حضور عامل در همان ناحیه را دارد. زیرا اگر عامل نسبت به محیط خود شناخت کامل‌تری داشته در هنگام تلاش برای رسیدن به اهداف خود به علت شناخت خوبی که از محیط دارد کمتر در محیط پرسه می‌زند (کمتر تلاش می‌کند) و با تعداد گام کمتری به سمت اهداف خود حرکت می‌کند - در واقع مسیر بهتری/کوتاه‌تری برای رسیدن به هدف را می‌شناسد. این موضوع در نهایت منجر می‌شود که عاملی که در هر ناحیه خبره‌تر است در همان ناحیه میزان پرسه زدن (حضور/تلاش) کمتری نسبت به دیگر عامل‌ها که از خبرگی نسبی کمتری برخوردار است را داشته باشد.

معیار تعریف شده در تعریف ۴-۱ قبلاً به صورت جزئی توسط احمدآبادی و همکاران [۱۲] ارائه شده است ولی معیار تعریف شده در این پژوهش تفاوت‌هایی با معیار احمدآبادی و همکاران دارد که به شرح زیر است:

۱. **میانگین تعداد قدم‌ها:** احمدآبادی و همکاران میانگین تعداد قدم‌های رسیدن به هدف (یا طبق تعریف ۴-۱ میانگین میزان ارجاع عامل در کل محیط - در زمانی که کل محیط را یک ناحیه در نظر بگیریم) را به عنوان معیار خبرگی در نظر گرفته‌اند در حالی که در تعریف ۴-۱ حرفی از میانگین آورده نشده است. ایرادی که معیار احمدآبادی و همکاران دارد این است که هنگامی که می‌خواهیم خبرگی عامل‌ها را بسنجیم صحیح نیست میانگین تعداد گام‌ها در نظر بگیریم زیرا ممکن است عامل در ابتدا بسیار نادان

بوده ولی بعد از طی مدتی به وسیله‌ی تجاربی خاص به عاملی بسیار دانا تبدیل شود و اگر میانگین‌گیری صورت گیرد آنگاه نادانی گذشته به میزان خبرگی کنونی تاثیر گذاشته و خبرگی عامل کمتر از میزان واقعی تخمین زده شود. در تعریف ۴-۱ خبرگی کنونی عامل مورد نظر است و کاری با مسیری که عامل برای کسب خبرگی کنونی‌اش طی کرده است نداریم.

۲. **انعطاف:** معیار احمدآبادی و همکاران از انعطاف برخوردار نیست و در خبرگی عامل‌ها را بصورت میانگین خبرگی در کل محیط محاسبه می‌کند در حالی که طبق تعریف ۴-۱ خبرگی عامل در نواحی مختلف از محیط قابل محاسبه است و همانطور که بعدها خواهیم دید خبرگی عامل‌ها در هر ناحیه به عنوان معیاری برای ترکیب دانش عامل‌ها نسبت به آن ناحیه مورد استفاده واقع خواهد؛ زیرا که عاملی ممکن است در حالت کلی محیط را آنچنان نشناخته باشد ولی در یک یا چند ناحیه بخصوص این عامل شناخت کامل‌تری از آن نواحی داشته باشد که معیار احمدآبادی و همکاران نمی‌تواند این مساله را در نظر بگیرد.

تا به اینجا گفته شد که عاملی که از خبرگی بیشتری برخوردار است لزوماً کمتر در محیط پرسه می‌زند و با طی کردن مسیر کوتاه‌تر به سمت اهداف خود، تلاش کمتری می‌کند ولی چند سوال در اینجا مطرح می‌شود که برای حل مساله نیازمند پاسخ به آن‌ها هستیم.

۱. میزان حضور عامل در نواحی مختلف، در محیطی که از  $d$ -بعد تشکیل شده است چگونه مدل شود؟
  ۲. اگر عاملی که در هر چرخه یادگیری به یکی از نواحی کلا وارد نشد و میزان پرسه زدن عامل در آن ناحیه صفر شود؛ آیا این مقدار کمینه پرسه زدن، نشان دهنده‌ی خبرگی عامل در آن ناحیه است؟
  ۳. چگونه در معیار خبرگی ارائه شده باید مساله عدم حضور عامل در یکی از نواحی را مدل کرد، بگونه‌ای که اثر سوئی بر تجربه‌ی دیگر عامل‌ها در آن نواحی، در هنگام ترکیب دانش عامل‌ها نداشته باشد؟
- پاسخ به این سوالات برای حل مساله با استفاده از معیار خبرگی پیشنهادی (تعریف ۴-۱) ضروری است. در پاسخ به سوال اول، ما به ازای کلیه‌ی نواحی یک ماتریسی به نام «ماتریس ارجاع» (یا به اختصار REFMAT<sup>۱</sup>) در نظر می‌گیریم که در ابتدا صفر مقداردهی شده‌اند و هر دفعه که عامل از حالتی به حالت دیگر می‌رود مقدار آن ناحیه‌ای که حالت جدید در آن واقع است یک واحد افزایش می‌یابد بدین وسیله میزان حضور عامل در نواحی مختلف را می‌شماریم. همانطور که در قسمت آزمایش‌های این پایان‌نامه نشان داده شده است میزان ریز یا درشت بودن این نواحی در کیفیت نتیجه تاثیرگذار نیست! یعنی عملاً چه ما در حالت کلی، کل محیط را به عنوان یک ناحیه در نظر بگیریم و میزان حضور عامل در این ناحیه را بشماریم (که معادل می‌شود با تعداد گام‌های عامل در طی رسیدن به هدف) یا در حالت جزئی به ازای هر حالت موجود را یک ناحیه در نظر بگیریم (که معادل می‌شود

<sup>1</sup>Reference Matrix

با تعداد ملاقات هر یکی از موقعیت‌ها توسط عامل) به یک نتیجه می‌رسیم. با توجه به دلیل فوق‌الذکر در پاسخ به سوال دوم، اگر تعداد نواحی زیاد باشد (مثلاً هر موقعیت یک ناحیه باشد - حداکثر تعداد نواحی) ممکن است عامل در طی رسیدن به هدف برخی از نواحی را کلاً ملاقات نکند و مقدار ارجاع به آن نواحی صفر شود و از طرفی طبق تعریف ۴-۱ عاملی که تعداد حضور کمتری در نواحی مختلف داشته باشد از خبرگی بیشتری در آن نواحی برخوردار است و در این شرایط که مقدار ارجاع عامل به ناحیه‌ای صفر باشد را نمی‌توان به خبرگی عامل در آن ناحیه نسبت داد زیرا که آن عامل در کل، آن ناحیه را ملاقات نکرده است که بخواهد تجربه‌ای را در تعامل با آن ناحیه کسب کند تا بتواند خبرگی خود را در آن ناحیه افزایش دهد. برای حل این مشکل و پاسخ به سوال سوم، ماتریسی جدیدی به نام ماتریس خاطره (یا به اختصار RCMAT<sup>۱</sup>) را معرفی می‌کنیم. این ماتریس وظیفه‌ی نگهداری آخرین تعداد ارجاعات عامل به هر کدام از نواحی تعریف شده را دارد و در زمان‌هایی که مقدار یک ناحیه در ماتریس REFMAT صفر باشد مقدار آن ناحیه از ماتریس RCMAT بروز رسانی می‌شود که میزان پرسه زدن عامل در آن ناحیه در آخرین باری عامل آن ناحیه را ملاقات کرده است را نشان می‌دهد؛ در صورتی که مقدار پرسه زدن یک ناحیه در ماتریس REFMAT مقداری غیر صفر باشد مقدار ماتریس RCMAT با مقدار کنونی REFMAT آن ناحیه بروز رسانی می‌شود.

دلیل استفاده از ماتریس RCMAT این است که در یادگیری تقویتی عامل زمانی می‌توان دانش (سیاست/خبرگی) خود را نسبت به نحوه‌ی عمل در یک موقعیت بهبود ببخشد که آن موقعیت را ملاقات کند. حال اگر عامل موقعیتی را ملاقات نکند دانش وی در آن موقعیت ثابت خواهد ماند به همین دلیل اگر عامل ناحیه‌ای را ملاقات نکند و مقدار REFMAT آن ناحیه صفر باشد می‌دانیم که دانش (خبرگی) عامل در آن ناحیه در این چرخه‌ی یادگیری ثابت مانده است و در صورتی که دوباره در آن ناحیه قرار می‌گرفت، حدوداً به همان میزان آخرین ملاقات در آن محیط پرسه خواهد زد. به عبارت دیگر در یک چرخه یادگیری ناحیه‌ی ملاقات نشده، تقریباً به میزان آخرین تعداد ارجاع شده برای آن ناحیه، مورد ارجاع واقع می‌شد.

#### ۴-۳ یادگیری مشارکتی $Q$ با استفاده از ماتریس ارجاع و انتگرال فازی

آنچه که تا به اکنون در مورد روش پیشنهادی این پژوهش آورده شده، معرفی یک معیار خبرگی است، که در تمامی موقعیت‌های دنیای واقعی به وفور مشاهده می‌شود و آن ارائه این فرضیه است که عامل خبره‌تر برای رسیدن به یک مجموعه از اهداف تلاش نسبی کمتری نسبت به دیگر عامل‌ها با خبرگی کمتر در شرایط یکسان می‌کند. حال که معیاری برای میزان خبرگی عامل‌ها در اختیار داریم چالش بعدی برای بهبود کیفیت و سرعت یادگیری مشارکتی ارائه‌ی روشی برای ترکیب دانش‌های عامل‌ها از محیط (جداول  $Q$  آن‌ها) با استفاده از معیار ارائه شده

<sup>۱</sup> Recall Matrix

می‌باشد. روش ترکیب باید بگونه‌ای باشد که کیفیت و سرعت یادگیری مشارکتی عامل‌ها را در طی زمان نسبت به زمانی که عامل‌ها بدون مشارکت یاد می‌گیرند بهتر کند. همچنین کیفیت و سرعت یادگیری همبستگی مستقیمی داشته باشند با تعداد عامل‌هایی که در حال اشتراک گذاری هستند؛ به عبارت دیگر در صورت افزایش تعداد عامل‌هایی که دانش‌های خود را به اشتراک می‌گذارند مدل ترکیب کننده‌ی دانش‌های آن عامل‌ها باید بتواند دانش بهتری تولید کند که نهایتاً منجر به بهتر شدن کیفیت و سرعت کلی یادگیری عامل‌ها شود.

در این پژوهش ما انتگرال فازی را به عنوان مدل ترکیب کننده‌ی دانش‌های عامل‌ها پیشنهاد می‌دهیم. دلیل انتخاب این مدل ویژگی‌های منحصر به فردی است که این مدل کننده در اختیار دارد که مدل را کاملاً مناسب برای ترکیب دانش عامل‌ها می‌کند. که در بخش‌های آتی فصل این ویژگی‌ها و دلایل مناسب بودن آن‌ها برای ترکیب دانش عامل‌ها آورده شده است. لازم به یادآوری است که همانطور که در قسمت ۳-۷ این پایان‌نامه آورده شده است ما به دلایل فنی از انتگرال فازی چوکت استفاده می‌کنیم که در بخش‌های بعدی این دلایل نیز بطور مفصل شرح داده می‌شود.

#### ۴-۳-۱ الگوریتم پیشنهادی

در این قسمت به معرفی الگوریتم پیشنهادی می‌پردازیم. آنچه که در الگوریتم ۴-۱ آمده است از دو قسمت تشکیل شده است، یک قسمت مربوط به یادگیری مستقل (خطوط ۸ تا ۱۳) و قسمت دیگر مربوط به یادگیری مشارکتی (خطوط ۱۵ تا ۲۳) می‌باشد. ورودی الگوریتم تعداد عامل‌ها می‌باشد و در ابتدا ماتریس‌های  $Q$  و REFMAT و RCMAT مقداردهی می‌شود. سپس تا زمانی که یادگیری پایان نیافته است ابتدا عامل‌ها در قسمت یادگیری مستقل به صورت جدا گانه در محیط فعالیت می‌کنند که رویه‌های آورده شده در خطوط ۸ تا ۱۲ همان الگوریتم یادگیری  $Q$  متعارف می‌باشد [۳۱]. در قسمت یادگیری مستقل تنها خط ۱۳ می‌باشد که در روش پیشنهادی به شبکه‌کد اضافه شده است و این تنها یک وظیفه‌ی بسیار ساده را انجام می‌دهد و آن شمارش میزان حضور عامل در هر کدام از نواحی از پیش تعیین شده است؛  $\phi(0)$  یک تابع نگاشت از یک موقعیت به یک ناحیه از محیط می‌باشد.

بعد از طی یادگیری مستقل عامل‌ها به قسمت اشتراک گذاری دانش‌های خود (جداول  $Q$ ) می‌رسند (خطوط ۱۵ تا ۲۳). در قسمت یادگیری مشترک ابتدا طبق آنچه که در در بخش قبلی آورده شده است جداول REFMAT و RCMAT به صورت مشترک بروزرسانی می‌شود و سپس جداول  $Q$  و REFMAT تمامی عامل‌ها به مدل ترکیب کننده فازی معرفی شده در این پژوهش فرستاده می‌شود و مدل ترکیب کننده فازی وظیفه‌ی استخراج یک دانش جدید با در نظر گرفتن ورودی‌های آن برای جایگزینی دانش قابلی عامل‌ها را دارد. لازم به ذکر است که الگوریتم ۴-۱ به صورت یک الگوریتم غیر-متمرکز<sup>۱</sup> بروی هر عامل می‌باشد [۲، ۱۲، ۱۹].

<sup>۱</sup>Decentralized

---

**الگوریتم ۴-۱** الگوریتم پیشنهادی یادگیری مشارکتی بر مبنای ماتریس REFMAT و انتگرال فازی
 

---

```

1: procedure REFMAT-COOPERATIVE-LEARNING( $m$ )
2:   Require:  $m > 1$  ▷ The number of agents.
3:   Ensure: Initialize the  $Q$  matrix;
4:   Ensure: Initialize the RCMAT  $\leftarrow 0$ ;
5:   Ensure: Initialize the REFMAT  $\leftarrow 0$ ;
6:   while not End Of Learning do
7:     if In individual learning mode then
8:       Visit the state  $s$ ;
9:       Select an action  $a$  based on an action selection policy;
10:      Carry out the  $a$  and observe a reward  $r$  at the new state  $s'$ ;
11:       $Q[s, a] \leftarrow Q[s, a] + \alpha(r + \lambda \max_{a'} (Q[s', a']) - Q[s, a])$ ;
12:       $s \leftarrow s'$ ;
13:      Increment REFMAT( $\phi(s)$ ) by one;
14:     else if In cooperative learning mode then
15:        $\vec{K} \leftarrow \{\}$ ;
16:        $\vec{R} \leftarrow \{\}$ ;
17:       for each agent  $i \leftarrow 1, m$  do
18:         REFMAT $_i$ , RCMAT $_i \leftarrow$  Conditional_Swap(REFMAT $_i$ , RCMAT $_i$ );
19:          $\vec{K}.add(Q_i)$ ;
20:          $\vec{R}.add(REFMAT_i)$ ;
21:       end for
22:        $Q \leftarrow$  FCI_Combiner( $\vec{K}$ ,  $\vec{R}$ );
23:       REFMAT  $\leftarrow 0$ ;
24:     end if
25:   end while
26: end procedure

```

---

الگوریتم تابع Conditional\_Swap(·) بسیار ساده می‌باشد و مقادیر غیر صفر ماتریس ارجاع را در ماتریس خاطره کپی می‌کند و مقادیر صفر ماتریس ارجاع را از ماتریس خاطره جایگزین می‌کند. این تابع در الگوریتم ۴-۲ آمده است.

در این پژوهش در دو قسمت نوآوری صورت گرفته است، قسمت اول ارائه‌ی معیاری جدید برای سنجش معیار خبرگی که طبق تعریف ۴-۱ این معیار در خط ۱۳ الگوریتم ۴-۱ پیاده‌سازی شده است؛ نوآوری دوم نحوه‌ی ترکیب اطلاعات دانش عامل‌ها با استفاده از انتگرال فازی که در خط ۲۲ الگوریتم ۴-۱ و شرح جزئیات پیاده‌سازی آن در الگوریتم ۴-۳ آمده است.

ورودی‌های الگوریتم ۴-۳ به ترتیب مجموعه‌ای از جداول  $Q$  و ماتریس‌های ارجاع (REFMAT) تمامی عامل‌ها می‌باشد بطوری که در ازای هر جدول  $Q$  یک ماتریس REFMAT متناظر وجود دارد. خروجی این الگوریتم یک جدول  $Q$  می‌باشد که از ترکیب جداول  $Q$  ورودی با در نظر گرفتن میزان خبرگی هر کدام از عامل‌ها که توسط ماتریس‌های REFMAT آن‌ها تعیین می‌شود، بدست آمده است. الگوریتم ۴-۳ به ازای کلیه‌ی



---

```

1: function Conditional_Swap(REFMAT, RCMAT)
2:   Require: size(REFMAT) = size(RCMAT)
3:   for each element  $r$  in REFMAT and its corresponding element  $c$  in RCMAT do
4:     if  $r = 0$  then
5:        $r = c$ ;
6:     else
7:        $c = r$ ;
8:     end if
9:   end for
10:  return REFMAT, RCMAT
11: end function

```

---

موقعیت‌ها ( $s$ ها در خط ۴) ابتدا مقادیر REFMAT کلیه‌ی عامل‌ها در ناحیه‌ای که آن موقعیت در آن واقع است (که توسط تابع نگاشت  $\phi(\cdot)$  بدست می‌آید) را استخراج می‌کند و در برداری بنام  $\vec{a}$  ذخیره می‌کند (خطوط ۶ و ۷) که در واقع میزان ارجاعات هرکدام از عامل‌ها در ناحیه‌ی  $\phi(s)$  می‌باشد.

بردار  $\vec{a}$  معیاری برای سنجش میزان خبرگی کلی عامل‌ها در موقعیت  $s$  است، طبق آنچه که در تعریف ۱-۴ آمده است در هر ناحیه عاملی خبره‌تر است که مقدار REFMAT مربوط به آن ناحیه از دیگر عامل‌ها کمتر باشد. از آنجایی که دامنه‌ی خروجی توابع  $g(\cdot)$  بازه‌ی  $[0, 1]$  می‌باشد و طبق آنچه که در الگوریتم‌های ۴-۴ تا ۷-۴ آورده شده است خروجی تابع  $g(\cdot)$  بر اساس ورودی‌های آن تخمین زده می‌شود و از طرفی نیاز به مکمل‌سازی میزان خبرگی عامل‌ها (عاملی که بیشترین ارجاع به ناحیه‌ای را دارد، کمترین میزان خبرگی را در آن ناحیه دارد و برعکس) نیاز به عادی‌سازی مقادیر ارجاعات داریم، در نتیجه در خط ۹ بعد از عادی‌سازی<sup>۱</sup> مقادیر REFMAT عامل‌ها در ناحیه‌ی  $\phi(s)$  یک مکمل‌گیری صورت می‌گیرد تا عاملی که مقدار REFMAT کمتری دارد دارای بیشترین مقدار بعد از عادی‌سازی شود. در خط ۱۰ به ازای کلیه‌ی عمل‌های ممکن در موقعیت  $s$  ابتدا مقادیر  $Q$  تک‌تک عامل‌ها را در موقعیت  $s$  و عمل  $a$  در خطوط ۱۲ و ۱۳ در بردار  $\vec{x}$  ذخیره می‌کنیم و در نهایت در خط ۱۵ با استفاده از انتگرال فازی چوکت معرفی شده در ۳-۹ مقدار  $Q$  مشارکتی حاصل از میزان خبرگی بردار  $\vec{A}$  و مقادیر  $Q$ های تک‌تک عامل‌ها در بردار  $\vec{x}$  در موقعیت  $s$  و عمل  $a$  بدست محاسبه می‌شود.

#### ۲-۳-۴ تعیین توابع $f(\cdot)$ و $g(\cdot)$ در انتگرال فازی چوکت

بطور خلاصه در الگوریتم ۳-۴ دو بخش عمده دارد بخش اول مربوط استخراج میزان خبرگی عامل‌ها بگونه‌ای که عاملی که خبره‌تر از دارای مقدار خبرگی بیشتری باشد که این بخش در خطوط ۶ تا ۹ صورت می‌گیرد؛ بخش دیگر محاسبه‌ی مقادیر  $Q$  مشارکتی کلیه‌ی عمل‌های ممکن در یک موقعیت با در نظر گرفتن میزان خبرگی عامل‌ها و مقادیر  $Q$  آن‌ها با استفاده از انتگرال فازی چوکت که در خطوط ۱۰ تا ۱۵ صورت می‌پذیرد.

---

<sup>1</sup>Normalize

---

```

1: function FCI_Combiner( $\vec{K}$ ,  $\vec{R}$ )
2:   Require:  $\text{length}(\vec{K}) = \text{length}(\vec{R}) = m$ 
3:   Ensure: Initialize  $\text{CoQ}_{\text{FCI}}$  ▷ The cooperative Q table.
4:   for each state  $s$  do
5:      $\vec{a} \leftarrow \{\}$ ; ▷ Contains the normalized values of REFMATs for state  $s$  for all agents.
6:     for each REFMAT $_i$  in  $\vec{R}$  do
7:        $\vec{a}.\text{add}(\text{REFMAT}_i(\phi(s)))$ ;
8:     end for
9:      $\vec{A} \leftarrow 1 - \text{normalize}(\vec{a})$ ;
10:    for each possible action  $a$  in state  $s$  do
11:       $\vec{x} \leftarrow \{\}$ ; ▷ Contains the  $Q$  values of action  $a$  in state  $s$  for all agents.
12:      for each  $Q_i$  in  $\vec{K}$  do
13:         $\vec{x}.\text{add}(Q_i[s, a])$ ;
14:      end for
15:       $\text{CoQ}_{\text{FCI}}[s, a] \leftarrow \sum_{i=1}^m (f(x_{\pi(i)}) - f(x_{\pi(i-1)})) \cdot g(\vec{A}_i)$  ▷ The Choquet Integral.
16:    end for
17:  end for
18:  return  $\text{CoQ}_{\text{FCI}}$ ;
19: end function

```

---

آنچه که در خط ۱۵ الگوریتم ۳-۴ مورد توجه واقع شود این است که توابع  $f(\cdot)$  و  $g(\cdot)$  چگونه باید تعریف شوند؟ برای تعیین تابع  $f(\cdot)$  منطقی که در این پژوهش استفاده کردیم بدین صورت است که از آنجایی که خروجی تابع  $g(\cdot)$  یک مقدار عددی<sup>۱</sup> بدون واحد می‌باشد و همچنین برای اینکه خروجی انتگرال فازی خط ۱۵ را بتوان به عنوان مقادیر جدول  $Q$  مشارکتی جدید در نظر گرفت تا بتوانیم در خطوط ۲۲ الگوریتم ۱-۴ به عنوان جدول  $Q$  تک‌تک عامل‌ها جایگذاری کنیم باید خروجی انتگرال فازی خط ۱۵ الگوریتم ۳-۴ از جنس جدول‌های  $Q$  عامل‌ها باشد در نتیجه تابع  $f(\cdot)$  باید یک تابع خطی بصورت ۱-۴ باشد تا خروجی انتگرال فازی همجنس مقادیر  $\vec{x}$  باشد.

$$f(\omega) = a\omega + b \quad (۱-۴)$$

متغیرهای  $a$  و  $b$  در ۱-۴ می‌تواند به عنوان پارامترهای وفقی<sup>۲</sup> در میزان کیفیت جدول  $Q$  مشارکتی خروجی الگوریتم ۳-۴ موثر واقع شود ولی با این حال در این پژوهش مقادیر  $a$  و  $b$  هر دو به ترتیب مقادیر ثابت ۱ و صفر در نظر گرفته شده‌اند که یعنی از تابع همانی به عنوان تابع  $f(\cdot)$  استفاده شده است؛ دلیل انتخاب مقادیر ۱ و صفر ثابت برای متغیرهای  $a$  و  $b$  این است که از بین تمام مقادیری که  $a$  و  $b$  می‌تواند بگیرند، استفاده از تابع همانی به

---

<sup>۱</sup>Scalar

<sup>۲</sup>Addaptive Parameters

عنوان تابع  $f(\omega)$  به نظر منطقی می‌رسد زیرا هیچ انتقال<sup>۱</sup> و درشت‌نمایی<sup>۲</sup> بروی مقادیر  $\omega$  (و نهایتاً مقادیر جدول  $Q$ ) اعمال نمی‌کند و به داده‌ها را به همان صورت که هستند می‌بیند – البته در صورت ارائه‌ی روشی که بتواند این مقادیر را به صورت وفقی تغییر دهد، شهود<sup>۳</sup> این را می‌گوید که می‌تواند در بهبود نتایج حاصله کمک کند؛ مثلاً در بازه‌های زمانی مشخص با تغییر مقدار درشت‌نمایی بتواند به افزایش سرعت یادگیری کمک کند.

تابع  $g(\cdot)$  یک ورودی مرتب شده طبق آنچه که در ۳-۱۱ آمده است می‌گیرد و در الگوریتم ۴-۳ تعیین این تابع تاثیر زیادی بروی کیفیت خروجی الگوریتم خواهد داشت ولی چالش‌هایی برای تعیین این تابع داریم؛ تابع  $g(\cdot)$  باید دارای ویژگی‌های زیر باشد:

۱. **پویا<sup>۴</sup> باشد:** از آنجایی که تابع  $g(\cdot)$  میزان اندازه‌گیری غیرافزایشی<sup>۵</sup> منابع اطلاعاتی را در اختیار می‌گذارد [۴]، نیاز داریم تعیین کنیم که کدام منابع اطلاعاتی (در اینجا خبرگی عامل‌ها) در کنار هم چه ارزش افزوده‌ای دارد؛ ولی از آنجایی که در حین یادگیری مشترک روشی برای تعیین این ارزش افزوده نداریم بنابراین باید تابع  $g(\cdot)$  بصورت پویا بتواند مقادیر این ارزش افزوده را تخمین بزند.

۲. **قابل گسترش<sup>۶</sup> باشد:** چون تعداد عامل‌ها در محیط متغیر است لذا باید تابع  $g(\cdot)$  بگونه‌ای باشد که به ازای تغییر تعداد عامل‌ها (که تغییر در تعداد اعضای بردار  $\vec{A}$  را در پی دارد) قابل گسترش باشد.

یکی از روش‌های تخمین  $g(\cdot)$  که دو ویژگی بالا را داشته باشد، تابع اندازه‌گیری- $\lambda$  سوگنو می‌باشد ولی این تابع نیاز به ریشه‌یابی روی متغیر  $\lambda$  دارد که طبق آنچه که در ۳-۵ آمده است به ازای تعداد عامل‌های مختلف نیاز به ریشه‌یابی معادلات غیرخطی دارد. بدلیل محاسبات سنگین و وقت‌گیر این ریشه‌یابی و همچنین نتایج حاصل از دستاوردهای این پژوهش (که در فصل نتیجه‌گیری آورده شده است)، در آزمایش‌ها صورت گرفته در این پژوهش از تابع اندازه‌گیری- $\lambda$  سوگنو به عنوان تابع  $g(\cdot)$  استفاده نشده است. یک سری توابع در این پژوهش بجهت استفاده، آزمایش و نتیجه‌گیری به عنوان  $g(\cdot)$  معرفی شده است که این توابع در الگوریتم‌های ۴-۴ تا ۴-۷ آمده‌اند.

در الگوریتم ۴-۴ به ازای هر ورودی دلخواه مقدار ثابت ۱ به عنوان خروجی برگشت داده می‌شود، این بدین معنی است که ارزش افزوده‌ی هر نوع ترکیب اطلاعاتی (خبرگی) برای ما دارای حداکثر ارزش می‌باشد و این مساله باعث می‌شود که نتیجه‌ی انتگرال فازی خط ۱۵ الگوریتم ۴-۳ مقداری معادل با مقدار خبره‌ترین عامل (عاملی که کمترین پرسه را در محیط مربوطه داشته) را به عنوان مقدار جدید جدول  $Q$  مشارکتی تولید کند.

<sup>1</sup> Shift

<sup>2</sup> Magnification

<sup>3</sup> Intuition

<sup>4</sup> Dynamic

<sup>5</sup> Non-additive

<sup>6</sup> Expandable

---

**الگوریتم ۴-۴** Const-One برای تخمین تابع  $g(\cdot)$  در الگوریتم ۳-۴
 

---

```

1: function Const-One( $\vec{A}_i$ )
2:   if  $\text{length}(\vec{A}_i) \geq m$  then
3:     return 1;
4:   else if  $\text{length}(\vec{A}_i) = 0$  then
5:     return 0;
6:   else
7:     return 1;
8:   end if
9: end function

```

---



---

**الگوریتم ۵-۴** Max برای تخمین تابع  $g(\cdot)$  در الگوریتم ۳-۴
 

---

```

1: function Max( $\vec{A}_i$ )
2:   if  $\text{length}(\vec{A}_i) \geq m$  then
3:     return 1;
4:   else if  $\text{length}(\vec{A}_i) = 0$  then
5:     return 0;
6:   else
7:     return  $\max_i(\vec{A}_i)$ ;
8:   end if
9: end function

```

---

در الگوریتم ۵-۴ میزان خبرگی خبره‌ترین عامل به عنوان خروجی تابع  $g(\cdot)$  برگشت داده می‌شود. در الگوریتم ۶-۴ خروجی، میانگین خبرگی عامل‌ها در نظر گرفته شده است و در الگوریتم ۷-۴ طبق رابطه‌ی نوشته شده میانگین  $k$ ام میزان خبرگی‌ها به عنوان خروجی برمی‌گردد به طوری که بزرگترین خبرگی در عدد  $k$  و کوچکترین خبرگی در عدد ۱ و هر آنچه که مابین این دو خبرگی وجود دارد در اندیس ترتیب مرتب شده آن‌ها ضرب می‌شود و میانگین این مجموع محاسبه و برگشت داده می‌شود.

نکته‌ای که در مورد الگوریتم ۶-۴ باید توجه کرد این است که با اینکه این الگوریتم شرط «یکنوایی» انتگرال چوکت را ارضا نمی‌کند ولی همانطور که در فصل آزمایش‌ها خواهیم دید باعث واگرایی الگوریتم پیشنهادی نمی‌شود که در این مساله جای تحقیق بیشتری دارد که شرایط لازم توابع  $g(\cdot)$  کاربرد مورد استفاده این پژوهش چگونه باید باشد، زیرا در الگوریتم ۶-۴ با نقض شدن شرط یکنوایی به همگرایی الگوریتم لطمه‌ای وارد نشد!

#### ۴-۴ علت کارکرد انتگرال فازی چوکت در انتقال دانش

در این قسمت به بررسی شهودی اینکه چرا انتگرال فازی چوکت برای انتقال (ترکیب) دانش‌های عامل‌ها می‌تواند موثر واقع باشد می‌پردازیم. این شهود بعدها در آزمایش‌ها نشان داده خواهد شد که صحت دارد. انتگرال فازی چوکت یک سری ویژگی‌ها دارد که برای مدل کردن انتقال دانش آن را کاندیدای مناسبی می‌کند. از مهم‌ترین

---

**الگوریتم ۴-۶** Mean برای تخمین تابع  $g(\cdot)$  در الگوریتم ۳-۴
 

---

```

1: function Mean( $\vec{A}_i$ )
2:   if  $\text{length}(\vec{A}_i) \geq m$  then
3:     return 1;
4:   else if  $\text{length}(\vec{A}_i) = 0$  then
5:     return 0;
6:   else
7:     return  $\frac{\sum_{j=1}^{\text{length}(\vec{A}_i)} \vec{A}_i(j)}{\text{length}(\vec{A}_i)}$ ;
8:   end if
9: end function

```

---



---

**الگوریتم ۴-۷** K-Mean برای تخمین تابع  $g(\cdot)$  در الگوریتم ۳-۴
 

---

```

1: function K-Mean( $\vec{A}_i$ )
2:   if  $\text{length}(\vec{A}_i) \geq m$  then
3:     return 1;
4:   else if  $\text{length}(\vec{A}_i) = 0$  then
5:     return 0;
6:   else if  $\text{length}(\vec{A}_i) = 1$  then
7:     return  $\vec{A}_i(1)$ ;
8:   else
9:      $\vec{B}_i = \text{Sort-Ascending}(\vec{A}_i)$ ;
10:    return  $\min\left\{\frac{\sum_{k=1}^{\text{length}(\vec{B}_i)} k \cdot \vec{B}_i(k)}{\left(\sum_{j=1}^{\text{length}(\vec{B}_i)} j\right) - 1}, 1\right\}$ ;
11:   end if
12: end function

```

---

ویژگی‌ها می‌توان به موارد زیر اشاره کرد [۲۶].

۱. **محدود است:** اگر شرایط مرزی تابع  $g(\cdot)$  برقرار باشد انتگرال فازی هیچ‌گاه بیشتر از حداکثر مقدار  $f(x_{\pi_i})$ ‌ها و کمتر از حداقل مقدار آن‌ها خروجی نمی‌دهد [۲۷]. یعنی دانش تولیدی خارج از محدوده‌ی دانش فعلی عامل‌ها نمی‌باشد فقط ترکیب مناسبی از این دانش‌ها به عنوان خروجی برگشت داده می‌شود که این در کاربرد یادگیری تقویتی به این معنی است که هیچ‌گاه مقادیر جدول  $Q$  بیشتر یا کمتر از آنچه که تجربه شده نمی‌شود؛ در نتیجه در صورت کران‌دار بودن پاداش‌های دریافتی از محیط جدول  $Q$  خروجی انتگرال فازی نیز کران‌دار است که نتیجه می‌دهد الگوریتم پیشنهادی حتماً همگرا خواهد شد [۱].

۲. می‌تواند اندازه‌گیری‌های غیرافزایشی را مدل کند: معمولاً روش‌هایی که تاکنون در این زمینه ارائه شده است از میانگین وزنی خبرگی عامل‌ها برای بدست آوردن جدول  $Q$  مشترک استفاده کرده‌اند [۲، ۱۲، ۱۹]. این در حالی هست که میانگین وزن‌دار قسمتی از مدل اندازه‌گیری غیرافزایشی می‌باشد. بنابراین با در نظر گرفتن مدل‌های غیرافزایشی که در ماهیت مساله هست قدرت و انعطاف بیشتری نسبت به روش‌هایی که

فقط از میانگین وزنی استفاده کرده‌اند، در اختیار داریم.

#### ۴-۴-۱ اثبات همگرایی روش پیشنهادی

در مورد همگرایی روش پیشنهادی می‌توان گفت که از آنجایی که در خطوط ۸ تا ۱۲ الگوریتم پیشنهادی ۴-۱ که یادگیری تقویتی  $Q$  بدون دخل و تصرفی آورده شده است، اثبات همگرایی یادگیری تقویتی عامل‌ها به قوت خود باقی است [۳۲]. حال باید نشان دهیم که ترکیب یادگیری تقویتی با انتگرال فازی همگرایی یادگیری تقویتی را برهم نمی‌زند.

*اثبات.* فرض می‌کنیم که به تعداد  $l$  عدد عامل وجود دارد که در هر مرحله‌ی یادگیری مشارکتی در زمان  $t$ ، مقادیر  $Q(s, a)$  عامل‌ها به ترتیب میزان بهینگی مقادیر آن‌ها به صورت زیر می‌باشند.

$$Q_{\pi_1}^t(s, a) \preceq Q_{\pi_2}^t(s, a) \preceq \dots \preceq Q_{\pi_l}^t(s, a) \preceq Q^*(s, a), \quad \forall s, a, t \quad (2-4)$$

که  $\pi_1$  اندیس عاملی است که نسبت به دیگر عامل‌ها دارای مقدار  $Q(s, a)$  با کمترین بهینگی می‌باشد و  $\pi_l$  اندیس بهینه‌ترین مقدار  $Q(s, a)$  است و عملگر  $\preceq$  به معنی «کمتر یا مساوی بودن از دیدگاه بهینگی» می‌باشد و  $Q^*(s, a)$  مقدار سیاست بهینه در  $(s, a)$  می‌باشد. طبق اثبات همگرایی یادگیری تقویتی (بدون در نظر گرفتن انتگرال فازی) زمانی که پاداش‌های دریافتی محیطی محدود و نرخ یادگیری محدود به بازه‌ی  $[0, 1]$  باشد آنگاه داریم [۳۲]:

$$Q_{\pi_i}^t(s, a) \rightarrow Q^*(s, a) \quad \text{as } t \rightarrow \infty, \quad \forall s, a, i \in [1 \dots l] \quad (3-4)$$

با توجه به اینکه که انتگرال فازی به ازای هر  $(s, a)$  مقداری مابین حداکثر و حداقل مقادیر جداول  $Q$  عامل‌ها در آن  $(s, a)$  را تولید می‌کند، طبق خاصیت محدود بودن انتگرال فازی چوکت می‌توانیم بگوییم که دانش تولیدی حاصل از ترکیب جداول  $Q_i, \forall i \in [1 \dots l]$  عامل‌ها در زمان  $t$  به جدول  $Q_{\text{REFMAT}}^t$  می‌رسیم که دارای خاصیت زیر است.

$$Q_{\pi_1}^t(s, a) \preceq Q_{\text{REFMAT}}^t(s, a) \preceq Q_{\pi_l}^t(s, a), \quad \forall s, a, t \quad (4-4)$$

از ۴-۲ تا ۴-۴ می‌توان به این نتیجه رسید از آنجایی که  $Q_{\pi_1}^t(s, a)$  به ازای هر  $t$  دلخواه و با توجه به

شرایط ذکر شده در اثبات همگرایی [۳۲] نهایتاً همگرا می‌شود و همچنین  $Q_{\pi_l}^t(s, a)$  نیز در  $t \rightarrow \infty$  همگرا می‌شود و کلیه‌ی جداول مابین این دو یعنی  $Q_{\pi_i}^t(s, a), \forall s, a, t, i \in [1 \dots l]$  نیز همگرا می‌شوند (۴-۲ و ۴-۳)؛ در نتیجه جدول تولیدی توسط انتگرال فازی چوکت که خاصیت ۴-۴ در آن برقرار می‌باشد نیز در  $t \rightarrow \infty$  همگرا خواهد شد.  $\square$

#### ۴-۵ نتیجه‌گیری

در این فصل به‌طور مفصل به معرفی روش پیشنهادی پرداختیم، ابتدا مفاهیم و تعاریف بنیادی روش پیشنهادی را ارائه دادیم، بعد از معرفی الگوریتم پیشنهادی به توضیح قسمت‌های مختلف این الگوریتم به‌صورت جز پرداختیم و در نهایت با اثبات همگرایی روش پیشنهادی به مبحث این فصل خاتمه بخشیدیم. در فصل بعدی آزمایش‌های انجام شده به‌جهت تأیید صحت روش پیشنهادی آورده شده است.

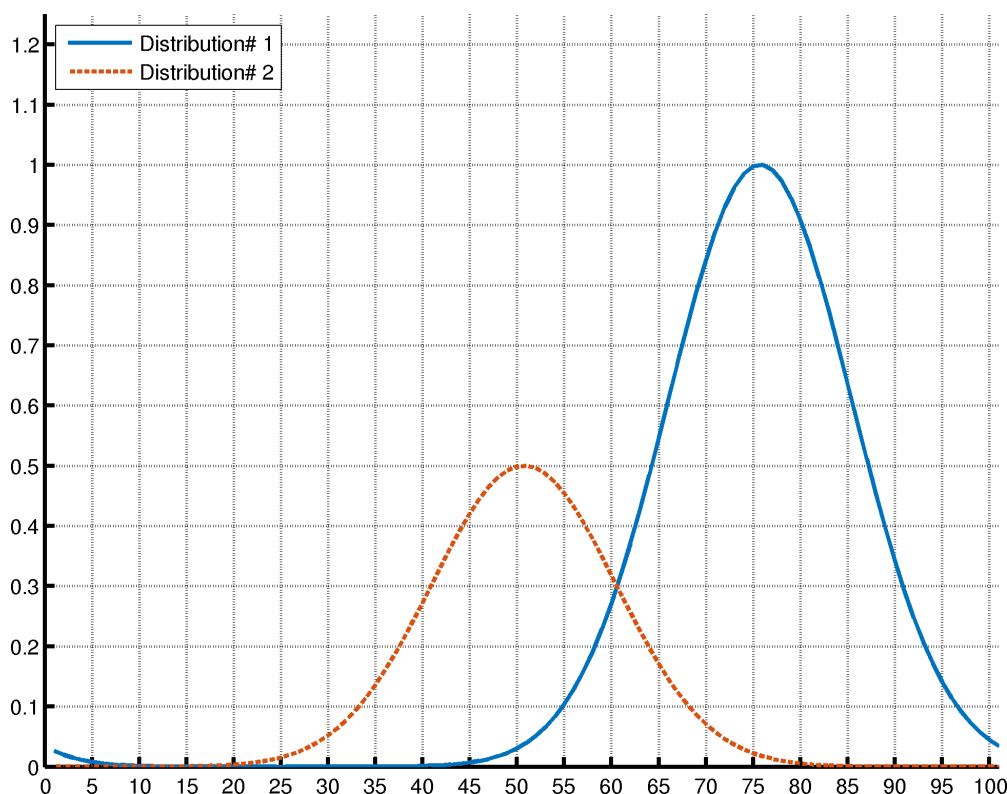
## فصل پنجم

### نتایج شبیه‌سازی و آزمایش‌ها

#### ۵-۱ مقدمه

در این فصل به ارائه‌ی آزمایش‌های صورت گرفته بروی روش پیشنهادی می‌پردازیم و در طی این آزمایش‌ها روش پیشنهادی را با روش کوتاه‌ترین مسیر تجربه شده (یا به اختصار SEP) مقایسه می‌کنیم که آخرین و مدرن‌ترین روش ارائه شده در جهت بهبود یادگیری مشارکتی می‌باشد [۲]. آزمایش‌ها بروی دو محیط «پلکان مارپیچ» و «صید و صیاد» صورت گرفته است. آزمایش‌ها به دو دسته تقسیم بندی شده است؛ دسته اول آزمایش‌هایی که روش پیشنهادی را در مقابل روش SEP قرار می‌دهد و عملکرد روش پیشنهادی را مورد سنجش قرار می‌دهد. دسته دوم آزمایش‌ها مربوط به آزمون رفتار روش پیشنهادی در صورت تغییر در پارامترهای مختلف آن می‌باشد. همچنین اثر استفاده از سیاست‌های انتخاب عمل مختلف در الگوریتم ۴-۱ نیز بررسی شده است. در روش‌های مرتبط قبلی [۲، ۱۹] که این پژوهش ادامه‌ی کار آن‌ها می‌باشد فقط از سیاست انتخاب عمل بولتزمن استفاده کرده‌اند؛ در این پژوهش علاوه بر بولتزمن تاثیر استفاده از روش  $\epsilon$ -حریصانه بروی هردو روش پیشنهادی و SEP نیز مورد بررسی واقع گردیده است.





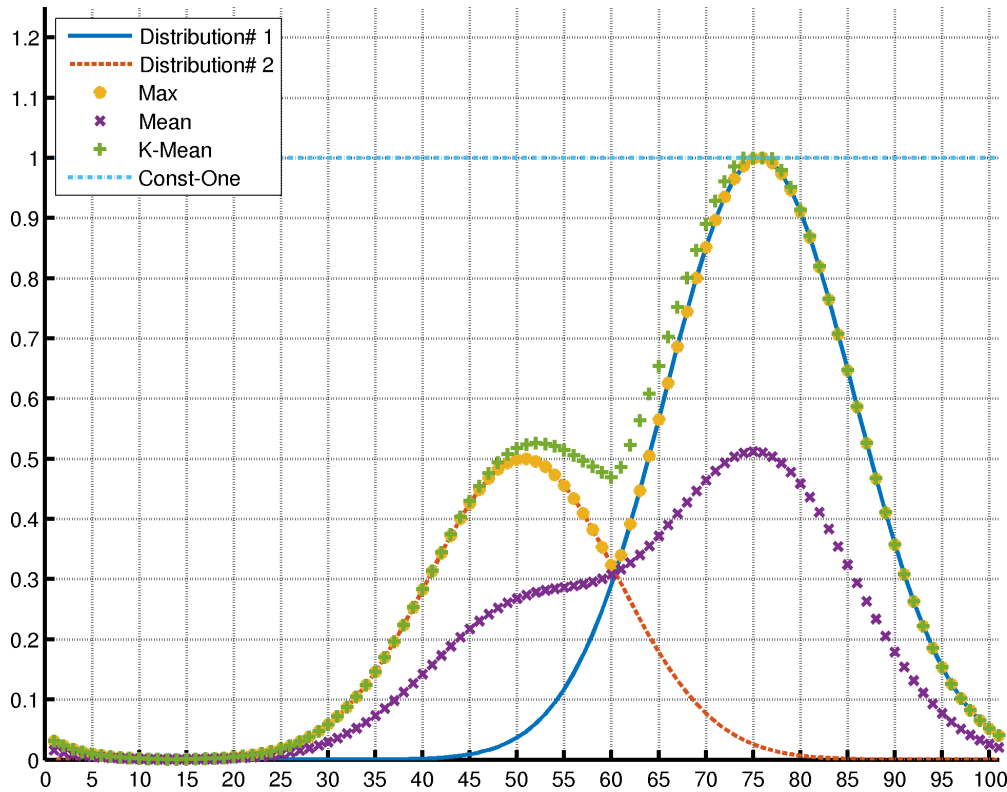
شکل ۵-۱: دو توزیع فرضی بجهت نمایش نحوه رفتار الگوریتم‌های ۴-۴ تا ۷-۴ بروی آن‌ها.

## ۵-۲ رفتار الگوریتم‌های معرفی شده برای $g(\cdot)$

در این قسمت به بررسی رفتار الگوریتم‌های ۴-۴ تا ۷-۴ معرفی شده برای  $g(\cdot)$  بروی دو توزیع فرضی خواهیم پرداخت، زیرا که در طی اجرای آزمایش‌های مختلف نتایج تاثیر این توابع بر اجرای الگوریتم پیشنهادی ۴-۱ آورده شده است، لذا به جهت درک علت تاثیرهای مختلف هر کدام از این توابع بروی نتیجه‌ی الگوریتم پیشنهادی در آزمایش‌ها، درک نحوه رفتار الگوریتم‌های ۴-۴ تا ۷-۴ ضروری است.

برای نمایش نحوه رفتار هر کدام از الگوریتم‌ها دو توزیع فرضی که در شکل ۵-۱ آورده شده است، را در نظر می‌گیریم. در صورت اعمال الگوریتم‌های ۴-۴ تا ۷-۴ بروی دو توزیع آورده شده در شکل ۵-۱ توزیع‌های جدیدی بصورت آنچه که در شکل ۵-۲ آمده است بدست می‌آیند. همانطور که در شکل ۵-۲ می‌بینیم اعمال الگوریتم Const-One بروی دو توزیع مقدار ثابت ۱ را برمی‌گرداند. اعمال الگوریتم Max در هر نقطه حداکثر مقدار هر دو توزیع را برمی‌گرداند. الگوریتم Mean میانگین دو توزیع را در هر نقطه حساب می‌کند و در نهایت الگوریتم K-Mean هر دو توزیع را محاسبه میکند که همانطور که می‌بینیم الگوریتم K-Mean به سبب ماهیت الگوریتم به سمت بیشترین مقدار پیش‌قدر<sup>۱</sup> می‌باشد.

<sup>۱</sup> Bias



شکل ۵-۲: نمایش توزیع‌های جدید بدست آمده بعد از اعمال الگوریتم‌های ۴-۴ تا ۷-۴ بروی دو توزیع فرضی شکل ۵-۱

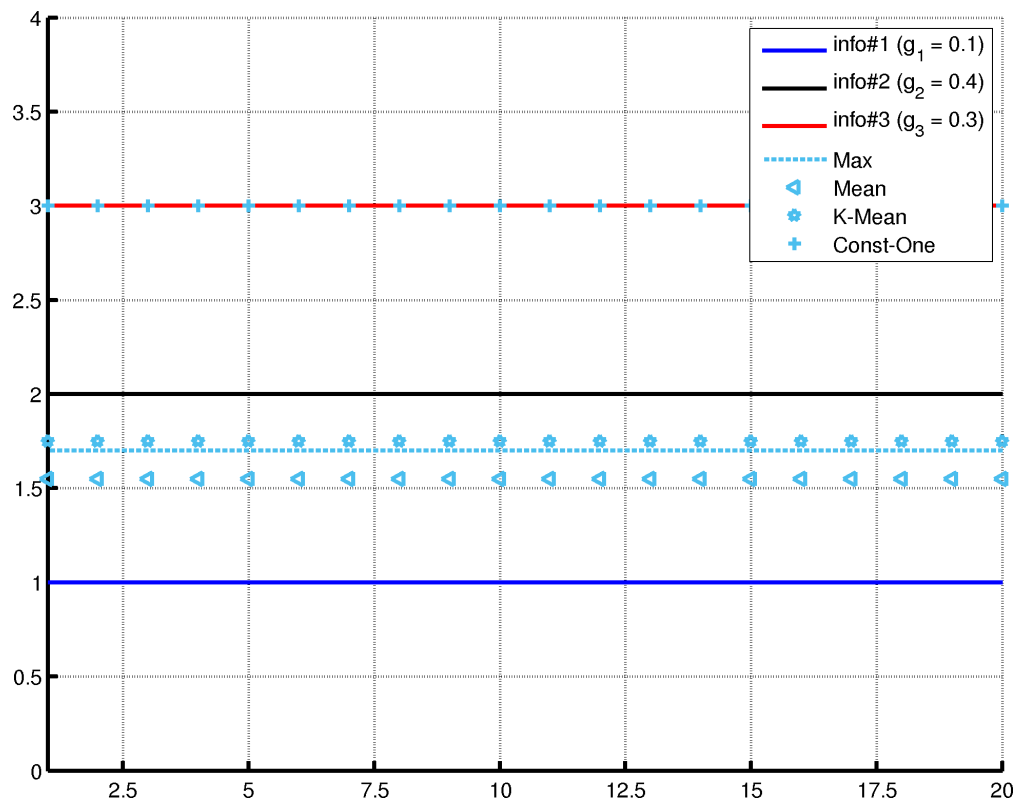
#### ۵-۲-۱ تعابیر مختلف انتگرال فازی چوکت از داده‌ها بر مبنای $g(\cdot)$

الگوریتم‌های ۴-۴ تا ۷-۴ به تنهایی فقط در نقش یک عملگر بازی می‌کند ولی در هنگام ترکیب دانش با انتگرال فازی چوکت به دانش خروجی الگوریتم از دیدگاه‌های متفاوتی نگاه می‌کنند. از آنجایی که در فصل‌های قبلی نیز آورده شد انتگرال فازی در واقع یک تعمیم الگوریتم دهنده‌ی میانگین وزنی می‌باشد که علاوه بر ویژگی‌هایی که روش میانگین وزنی ارائه می‌دهد می‌تواند اندازه‌گیری‌های غیرافزایشی را نیز مدل کند. لذا با تغییر تابع  $g(\cdot)$  می‌توان باعث شد که انتگرال فازی چوکت تعابیر مختلفی از داده‌های ورودی خود ارائه دهد. از بین الگوریتم‌ها فقط الگوریتم Const-One دارای تعبیر صریح ریاضی می‌باشد که در رابطه‌ی ۵-۱ آمده است، بقیه‌ی الگوریتم‌ها دارای تعابیر صریح نیستند و فقط می‌توانیم بر اساسی نمایشی که در شکل ۵-۲ آمده است شهودی از نحوه‌ی تغییر رفتار انتگرال فازی به ازای هریک از الگوریتم‌ها ارائه داد.

$$g = \text{Const-One}(\cdot) \equiv \begin{cases} g(X) & = 1 \\ g(\emptyset) & = 0 \\ g_{A \subseteq X}(A) & = 1 \end{cases} \Rightarrow C_g(f) \equiv \max\{f(x_{\pi(1)}^c), \dots, f(x_{\pi(n)}^c)\} \quad (۵-۱)$$

برای نمایش شهودی نحوه‌ی تغییر رفتار انتگرال فازی چوکت در شکل ۵-۳ سه منبع اطلاعاتی با مقادیر

$$g = [0.1 \quad 0.4 \quad 0.3]^T \text{ و } y = 2 \text{ و } y = 3 \text{ در نظر گرفته شده است و مقدار ارزش هر کدام از این‌ها به ترتیب } g = [0.1 \quad 0.4 \quad 0.3]^T$$



شکل ۳-۵: نمایش رفتار انتگرال فازی بروی منابع اطلاعاتی  $y = 1$  و  $y = 2$  و  $y = 3$  به ازای توابع  $g(\cdot)$ های مختلف.

در نظر گرفته شده است. سپس انتگرال فازی چوکت را با در نظر گرفتن تابع همانی به عنوان تابع  $f(\cdot)$  بر روی این ۳ منبع اطلاعاتی اعمال کردیم و همانطور که می‌بینیم مقداری که انتگرال فازی چوکت به ازای  $g = \text{Const-One}(\cdot)$  تولید می‌کند برابر با حداکثر مقدار منابع اطلاعاتی دریافتی می‌باشد. در حالت کلی هرچقدر میانگین تابع  $g_{A \subseteq X}(A)$  به سمت مقدار ۱ متمایل باشد خروجی انتگرال فازی چوکت به سمت بیشینه مقدار منابع اطلاعاتی پیش‌قدر می‌شود و در صورتی که این میانگین به سمت صفر متمایل باشد خروجی به کمینه مقدار پیش‌قدر می‌شود.

### ۳-۵ مقایسه‌ی روش پیشنهادی با روش کوتاه‌ترین مسیر تجربه شده

در این قسمت به مقایسه‌ی روش پیشنهادی با روش «کوتاه‌ترین مسیر تجربه شده» که بروزترین تکنیک ارائه شده در این شاخه از یادگیری مشارکتی می‌باشد می‌پردازیم [۲]. کلیه‌ی این آزمایش‌ها در دو محیط «پلکان مارپیچ» و «صید و صیاد» صورت گرفته است. نتیجه‌ی هر آزمایش حاصل میانگین ۲۰ اجرای مستقل تمامی الگوریتم‌ها می‌باشد. همچنین به غیر از مواردی که صراحتاً قید شده است تعداد عوامل ۳ عدد می‌باشد - البته بدیهی است که یادگیری مستقل تک عامله (یا به اختصار IL<sup>۱</sup>) شامل این قاعده نمی‌باشد. همچنین در کلیه‌ی

<sup>۱</sup>Individual Learning

جدول ۵-۱: لیست اختصارهای استفاده شده در این فصل

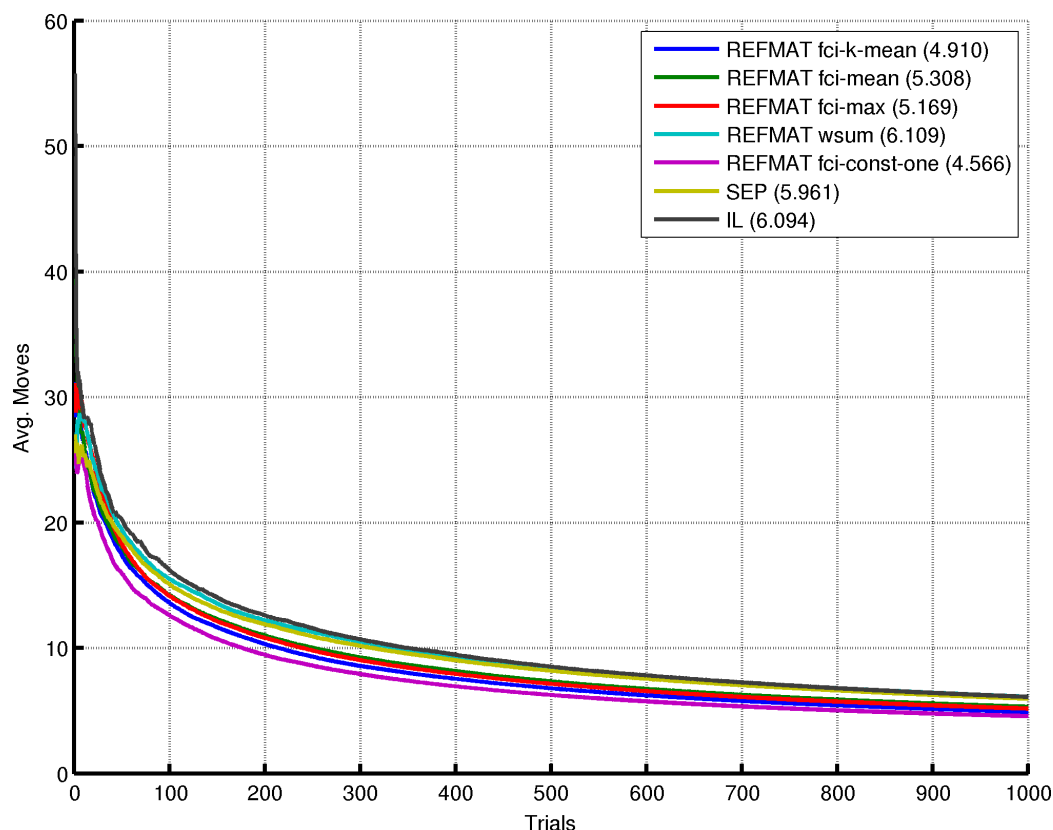
معنی	اختصار
روش پیشنهادی	REFMAT
یادگیری مستقل تک عامله	IL
روش کوتاه‌ترین مسیر تجربه شده	SEP
میانگین وزنی	wsum
الگوریتم Max به عنوان مدل کننده‌ی تابع $g(\cdot)$	fci-max
الگوریتم Mean به عنوان مدل کننده‌ی تابع $g(\cdot)$	fci-mean
الگوریتم K-Mean به عنوان مدل کننده‌ی تابع $g(\cdot)$	fci-k-mean
الگوریتم Const-One به عنوان مدل کننده‌ی تابع $g(\cdot)$	fci-const-one
جستجوی کاملاً تصادفی محیط	Rand-Walk

آزمایش‌ها عامل‌ها از ۲۰۰ چرخه یادگیری بهره می‌برند و در هر چرخه عامل ۵ بار تلاش<sup>۱</sup> مستقل می‌کند که در مجموع ۱۰۰۰ تلاش صورت می‌گیرد. کلیه‌ی پارامترهای مربوط قسمت یادگیری مستقل الگوریتم ۴-۱ اعمال شده در آزمایش‌ها این فصل منطبق بر پارامترهای تعریف شده کار میرزایی می‌باشد که نتایج قایل قیاس باشند [۲]. در ضمن در این فصل اختصارهای جدول ۵-۱ را نیز داریم.

در این فصل در حالت کلی ما در دو بخش سیاست انتخاب عمل «بولتزمن» و « $\epsilon$ -حریصانه» (که از این به بعد، به اختصار «تابع بولتزمن» و «تابع  $\epsilon$ -حریصانه» خطاب خواهیم کرد.) به مقایسه‌ی نتایج می‌پردازیم. طبق آنچه که در ادامه مشاهده خواهیم کرد چه در صورت استفاده از تابع بولتزمن و چه تابع حریصانه روش پیشنهادی چه در سرعت یادگیری و چه در کیفیت یادگیری بهتر از روش SEP می‌باشد.

برای اینکه نشان دهیم که استفاده از انتگرال فازی در بهبود نتیجه تاثیر بسزایی دارد از تابع میانگین وزنی (یا به اختصار wsum<sup>۲</sup>) نیز استفاده کرده‌ایم. بدین صورت که بجای اینکه بعد از استخراج میزان خبرگی هر عامل، جدول  $Q$  هر عامل را متناسب با میزان خبرگی‌ای که دارد در دانش جمعی دخیل می‌کنیم تا جدول  $Q$  مشارکتی تولید شود. تابع میانگین وزنی روشی است که در پژوهش‌های اخیر به کرات از آن استفاده کرده‌اند [۲، ۱۲، ۱۹]. یکی از اهداف ما در این پژوهش نمایش قدرت انتگرال‌های فازی می‌باشد.

<sup>۱</sup>Trial<sup>۲</sup>Weighted Sum



شکل ۴-۵: مقایسه در سرعت و کیفیت یادگیری با تابع بولتزمن در محیط پلکان مارپیچ

### ۵-۳-۱ مقایسه در محیط پلکان مارپیچ

آزمایش‌های مربوط به این قسمت در ۴ بخش صورت گرفته است؛ ۱. مقایسه در سرعت و کیفیت یادگیری، ۲. مقایسه در پیچیدگی زمانی، ۳. مقایسه در میزان باروری، ۴. مقایسه تاثیر تعداد عامل‌ها بر میزان کیفیت و سرعت یادگیری.

#### سیاست انتخاب عمل «بولتزمن»

مقایسه در سرعت و کیفیت یادگیری: نتایج حاصل از اجرای الگوریتم‌ها در محیط پلکان مارپیچ در شکل ۴-۵ آمده است. در این شکل محور افقی تعداد تلاش‌های یادگیری عامل را نشان می‌دهد که در تلاش اول عامل بدون دانش اولیه شروع به تعامل با محیط می‌کند و در تلاش ۱۰۰۰ام عامل به اجرای خود پایان می‌دهد. محور عمودی نمودار میانگین تعداد قدم‌های عامل را نشان می‌دهد. اعداد کناری برجسب‌ها (گوشه بالا سمت راست) متوسط تعداد قدم در آخرین تلاش عامل می‌باشد که انتظار می‌رود عامل آگاهی نسبی کاملی از محیط دارد را نشان می‌دهد که این عدد هرچقدر کمتر باشد نشان می‌دهد که عامل در طی رسیدن به هدف تعداد گام کمتری برداشته است و در نتیجه دانش و شناخت بهتری از محیط دارد.

همانطور که مشاهده می‌شود روش SEP دارای ۲٪ بهبود نسبت به IL می‌باشد در حالی که روش پیشنهادی

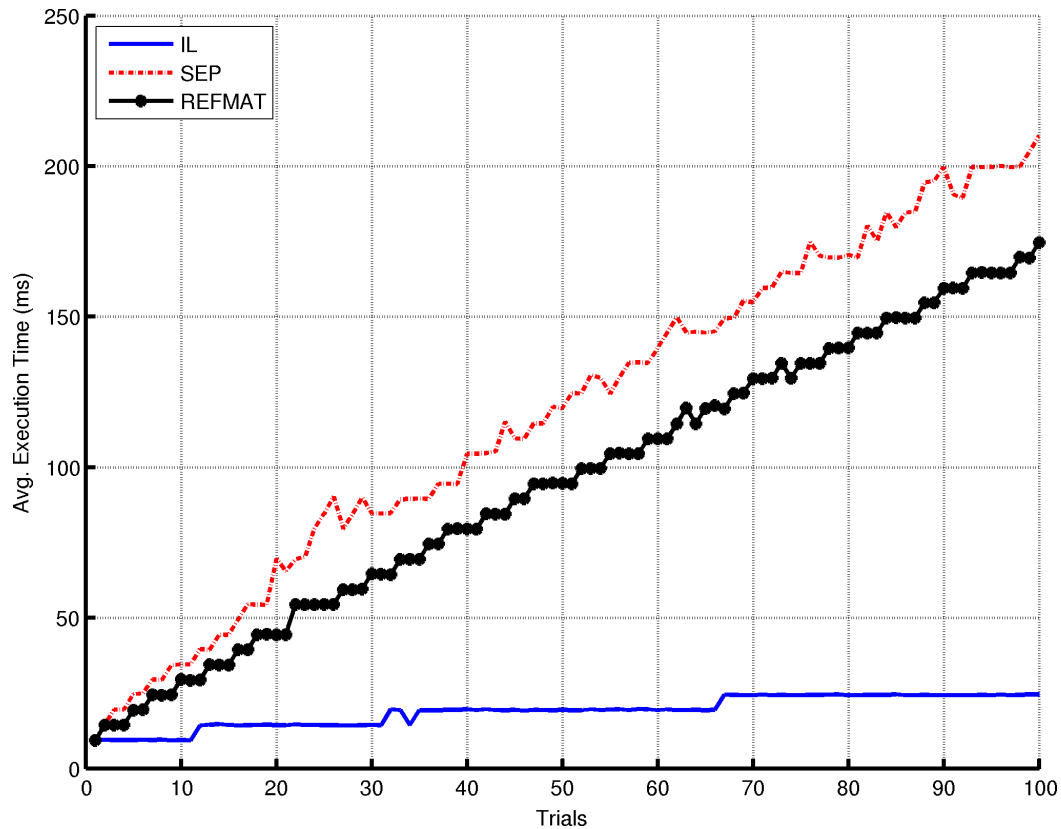
جدول ۵-۲: مقایسه در میزان درصد بهبود کیفیت یادگیری در محیط پلکان مارپیچ با تابع بولتزمن

			REFMAT				
	IL	SEP	wsum	fci-mean	fci-max	fci-k-mean	fci-const-one
IL	0.0						
SEP	2.2	0.0					
wsum	-0.2	-2.3	0.0				
fci-mean	14.9	12.5	15.1	0.0			
fci-max	18.0	15.5	18.2	2.7	0.0		
fci-k-mean	24.0	21.4	24.2	7.9	5.1	0.0	
fci-const-one	33.6	30.7	33.8	16.2	13.2	7.7	0.0

در زمانی که از انتگرال فازی استفاده می‌کند در بدترین حالت دارای ۱۸٪ بهبود و در بهترین حالات دارای ۳۳٪ بهبود می‌باشد که نسبت به روش SEP تقریباً ۹ الی ۱۶ برابر نتیجه را بهبود داده است. در صورتی که از میانگین وزنی بجای انتگرال فازی استفاده شود نتایج با اختلاف اندکی (کمتر از ۱-٪) بدتر از یادگیری IL بوده است که نشان می‌دهد که استفاده از انتگرال فازی چقدر می‌تواند نسبت به روش‌های سنتی و معمولی چون میانگین وزنی موثر واقع شود. نتایج این قسمت را می‌توان در جدول ۵-۲ خلاصه کرد. همچنین طبق آنچه در فصل ۳ آورده شده سرعت یادگیری به صورت مساحت زیر نمودار در شکل ۵-۴ محاسبه می‌شود، به صورت دیداری<sup>۱</sup> می‌توانیم ببینیم روش پیشنهادی دارای کمترین مساحت زیر منحنی می‌باشد که نشان می‌دهد دارای بیشترین سرعت یادگیری می‌باشد.

**مقایسه در سرعت اجرا:** در این قسمت به مقایسه‌ی پیچیدگی زمانی روش پیشنهادی با روش SEP مورد بررسی قرار می‌گیرد، برای محاسبه‌ی پیچیدگی زمانی به روش ریاضی کار بسیار دشوار و پرخطایی می‌باشد - زیرا محاسبه‌ی پیچیدگی تک‌تک بخش‌های الگوریتم که خودشان از زیربخش‌های مختلف و پیچیده‌ای تشکیل شده است کاری پر خطا می‌باشد؛ در اینجا ما بجای محاسبه‌ی پیچیدگی زمانی ریاضی دو الگوریتم از مدت زمانی که طول می‌کشد برنامه در سیستم اجرا و خاتمه یابد استفاده می‌کنیم. در شکل ۵-۵ میانگین زمانی ۲۰ اجرای مستقل برحسب میلی‌ثانیه به ازای هریک از تعداد تلاش‌ها آورده شده است. همان‌طور که در شکل ۵-۵ مشاهده می‌شود الگوریتم IL دارای حداکثر سرعت اجرا می‌باشد زیرا که هیچ سربار محاسباتی یادگیری مشترک را ندارد؛ هدف یادگیری اشتراکی این است که می‌خواهد در ازای یک سری سربار محاسباتی کیفیت و سرعت «یادگیری»

<sup>۱</sup> Visual



شکل ۵-۵: مقایسه در سرعت اجرای روش‌ها به ازای تعداد تلاش‌های متفاوت برحسب میلی‌ثانیه با تابع بولتزمن در محیط پلکان مارپیچ

عامل‌ها را افزایش دهد. با در نظر داشتن این موضوع همانطور که قبلاً دیدیم روش پیشنهادی سرعت و کیفیت یادگیری را بیشتر از روش SEP افزایش می‌دهد و در اینجا نیز می‌بینیم که دارای سرعت اجرای بیشتری نسبت به روش SEP می‌باشد که نشان از بهینه‌گی روش پیشنهادی نسبت به روش SEP می‌دهد.

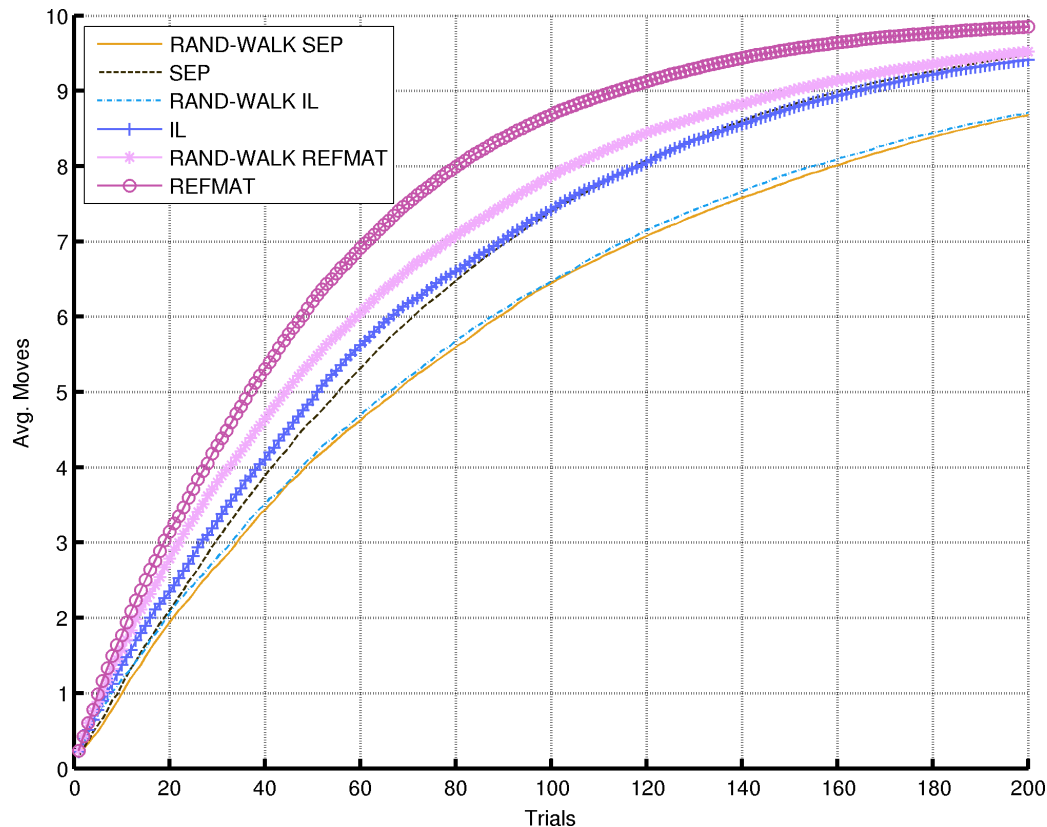
#### مقایسه در میزان باروری:

**تعریف ۵-۱ (باروری).** اگر فرض کنیم الگوریتم یادگیری تقویتی  $\psi_Q(\mathcal{E})$  وجود دارد که در محیط  $\mathcal{E}$  فعالیت می‌کند و دانش خود را در جدولی مانند  $Q$  ذخیره می‌کند، باروری الگوریتم  $\psi_Q(\mathcal{E})$  در هر لحظه را حداکثر مقدار جدول  $Q$  در آن لحظه تعریف می‌کنیم.

**تعریف ۵-۲ (سرعت باروری).** اگر فرض کنیم الگوریتم یادگیری تقویتی  $\psi_Q(\mathcal{E})$  وجود دارد که در محیط  $\mathcal{E}$  فعالیت می‌کند و دانش خود را در جدولی مانند  $Q$  ذخیره می‌کند، سرعت باروری الگوریتم  $\psi_Q(\mathcal{E})$  در هر لحظه را سرعت همگرایی حداکثر مقدار جدول  $Q$  به سمت حداکثر پاداش محیط قابل دریافت تعریف می‌کنیم.

**تعریف ۵-۳ (میزان باروری).** انتگرال سرعت باروری را میزان باروری الگوریتم  $\psi_Q(\mathcal{E})$  که در محیط  $\mathcal{E}$  فعالیت می‌کند و دانش خود را در جدولی مانند  $Q$  ذخیره می‌کند، تعریف می‌کنیم.

**فرضیه ۵-۱ (معیاری جدید برای سرعت یادگیری).** طبق تعاریف ۵-۲ و ۵-۳ الگوریتمی میزان باروری بیشتری دارد که



شکل ۵-۶: نمودار باروری الگوریتم‌ها مختلف با تابع بولتزمن در محیط پلکان مارپیچ

سرعت‌تر مقادیر جدول  $Q$  خود را به سمت بیشه مقداری که می‌تواند داشته باشد سوق دهد. معمولاً در الگوریتم‌های یادگیری تقویتی  $Q$  این کار با تنظیم مقدار سرعت یادگیری  $\alpha$  صورت می‌گیرد که باعث می‌شود الگوریتم‌ها با سرعت بیشتری به یادگیری نحوه تعامل با محیط بپردازند. لذا در شرایط یکسان می‌توان گفت الگوریتمی بهتر عمل می‌کند که نحوه تعامل با محیط را سریع‌تر نسبت به دیگر الگوریتم‌ها یاد می‌گیرد و میزان باروری بیشتری داشته باشد.

در شکل ۵-۶ آورده شده است حداکثر میزان جدول  $Q$  روش‌ها در هر تلاش آورده شده است. همانطور که قبلاً در تعریف محیط پلکان مارپیچ آورده شده است حداکثر مقدار پاداش این محیط مقدار ۱۰ می‌باشد لذا همان‌طور که مشاهده می‌شود الگوریتم‌ها با مساحت‌های زیر نمودار متفاوتی مقادیر جدول  $Q$  خود را به سمت حداکثر مقدار پاداش قابل دریافت از محیط سوق می‌دهد که نشان از میزان باروری الگوریتم‌ها می‌دهد.

در شکل ۵-۶ منظور از RAND-WALK حرکت کاملاً تصادفی می‌باشد، به این صورت که عامل بعد از هر حرکت جدول  $Q$  خود را بروز رسانی می‌کند ولی هنگام انتخاب عمل در تابع بولتزمن مقدار  $\tau \rightarrow +\infty$  در نظر گرفته می‌شود تا میزان احتمال تمامی حرکت‌ها یکسان شود و در نتیجه حرکتی به صورت تصادفی انتخاب شود. همان‌طور که در قسمت‌های قبل دیدیم روش پیشنهادی هم در کیفیت و هم در سرعت یادگیری بهبود چشم‌گیری دارد و از طرفی هم در نمودار ۵-۶ دارای بیشترین میزان باروری (مساحت زیر نمودار) حداکثر مقدار جدول  $Q$  می‌باشد که این مساله تایید کننده فرضیه ۵-۱ می‌باشد.

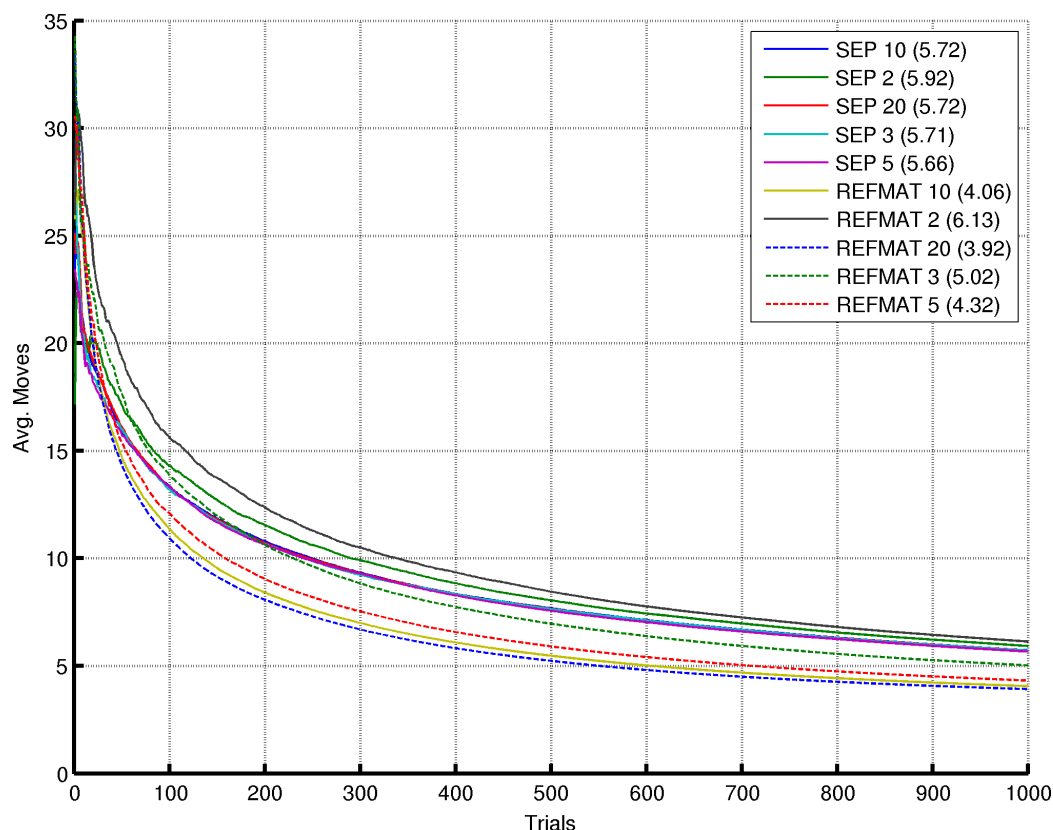


دلیل وجود نتایج آزمایش اجرای RAND-WALK در این قسمت این است که بررسی کنیم در صورتی که عامل بصورت کورکورانه حرکت کند روش معرفی شده و SEP چقدر در میزان بارور شدن جدول  $Q$  عامل‌ها موثرند؟ به عبارت دیگر، در صورتی که استراتژی خاصی جهت انتخاب عمل وجود نداشته باشد، روش‌ها چقدر قدرت باروری دارند؟ همانطور که در شکل ۵-۶ مشاهده می‌کنیم روش معرفی شده در زمانی که به صورت تصادفی اقدام به انتخاب عمل می‌کند بیشتر از زمانی که IL با استفاده از تابع بولتزمن اقدام به انتخاب عمل می‌کند جدول  $Q$  را بارور می‌کند که از قدرت روش ارائه شده خبر می‌دهد. همچنین در مورد روش SEP می‌بینیم که در زمانی که بصورت تصادفی اقدام به عمل می‌کند باروری کمتری نسبت به روش پیشنهادی و IL دارد؛ یعنی میزان باروری روش SEP وابستگی زیادی به سیاست انتخاب عمل دارد و در صورت نداشتن سیاست انتخاب عمل خاصی بشدت عملکردش کاسته می‌شود ولی در روش پیشنهادی میزان این وابستگی از شدت کمتری برخوردار است که از دیگر امتیازها مثبت روش پیشنهادی می‌باشد.

**مقایسه تاثیر تعداد عامل‌ها بر میزان کیفیت و سرعت یادگیری:** در این مقایسه سعی شده است که تاثیر یک فاکتور بنیادی سیستم‌های چندعامله مشارکتی را مورد بررسی قرار دهیم، و آن میزان تاثیر پذیری روش‌های مورد مقایسه با افزایش تعداد عامل‌ها می‌باشد. در تئوری سیستم‌های چندعامله مشارکتی دیدگاه معقول براین است که اثر تعداد عامل‌ها در کیفیت و سرعت یادگیری مشارکتی باید مثبت باشد. در غیر این صورت سیستم‌های چندعامله‌ای که تعداد عامل‌ها تاثیری در خروجی سیستم نداشته باشد، دیگر ماهیت سیستم‌های چندعامله را ندارد.

همان‌طور که در شکل ۵-۷ آمده است، روش پیشنهادی و روش SEP به ازای تعداد عامل‌های ۲، ۳، ۵، ۱۰ و ۲۰ عدد به تعداد ۲۰ بار اجرا درآمده و میانگین اجراها به نمودار کشیده شده است. همانطور که می‌بینیم روش SEP در زمانی که ۲۰ عامل در حال یادگیری و اشتراک گذاری دانش‌های خود هستند نسبت به زمانی که فقط ۲ عامل در حال تعامل مشارکتی با محیط هستند فقط ۳٪ در خروجی الگوریتم تاثیر مثبت داشته است. این در حالی است که در همین شرایط میزان بهبود نتیجه‌ی روش پیشنهادی ۵۶٪ می‌باشد. که نشان می‌دهد روش SEP نسبت به افزایش تعداد عامل‌ها رفتاری تقریباً خنثی از خود نشان می‌دهد درحالی که روش پیشنهادی در ازای افزایش تعداد عامل‌ها به دلیل اینکه دانش جمعی نیز افزایش می‌یابد کیفیت خروجی آن نیز بهتر می‌شود.

**نتیجه‌گیری:** نتیجه‌ای که از مقایسه‌ی روش پیشنهادی در هر چهار مقایسه‌ی بالا می‌توان گرفت این است که روش پیشنهادی بهبود چشم‌گیری به روش SEP در محیط پلکان مارپیچ و سیاست انتخاب عمل بولتزمن داده است.



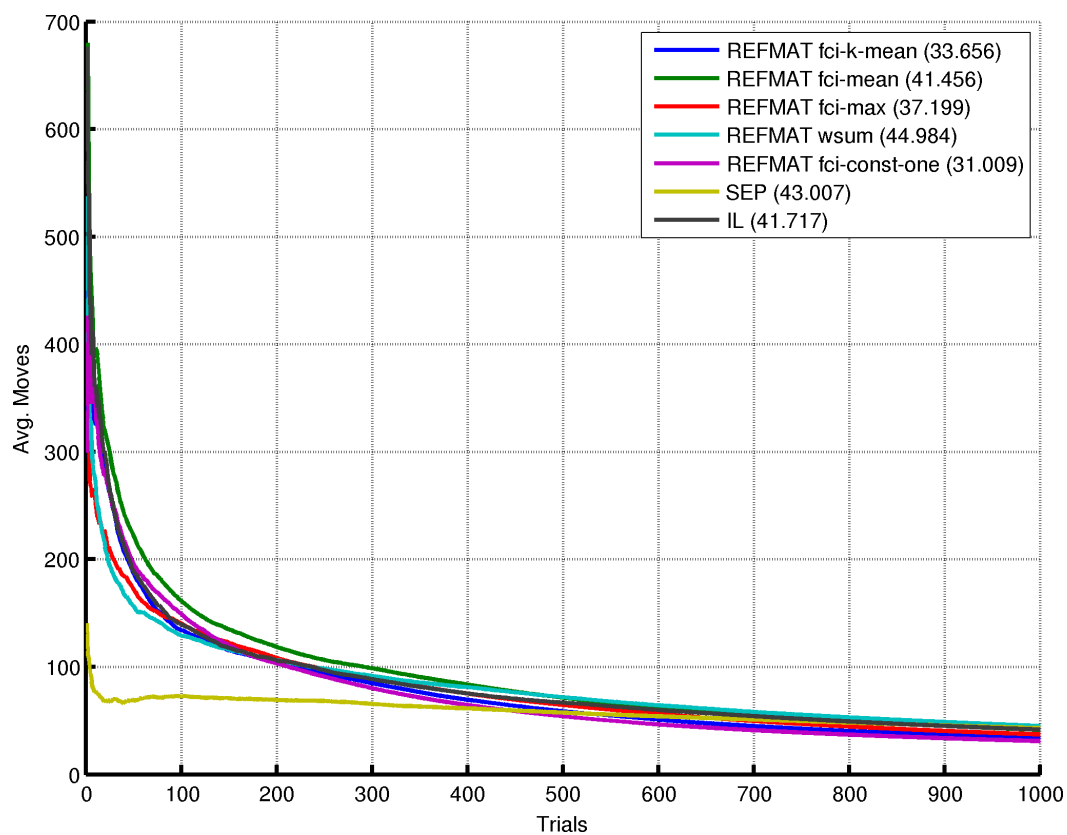
شکل ۵-۷: مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری با تابع بولتزمن در محیط پلکان مارپیچ

#### سیاست انتخاب عمل « $\epsilon$ -حریصانه»

مقایسه در سرعت و کیفیت یادگیری: نتایج حاصل از اجرای الگوریتم‌ها در محیط پلکان مارپیچ در شکل ۵-۸ آمده است. شرایط این آزمایش مشابه با شرایط آزمایش با تابع بولتزمن می‌باشد.

همانطور که مشاهده می‌شود روش SEP دارای ۳-٪ بهبود نسبت به IL می‌باشد در حالی که روش پیشنهادی در زمانی که از انتگرال فازی استفاده می‌کند در بدترین حالت دارای ۶/۰٪ بهبود و در بهترین حالات دارای ۳۴٪ بهبود می‌باشد که نسبت به روش SEP تقریباً ۴ الی ۳۸ برابر نتیجه را بهبود داده است. در صورتی که از میانگین وزنی بجای انتگرال فازی استفاده شود نتایج با اختلافی حدود ۷-٪ بدتر از یادگیری IL بوده است که نشان می‌دهد که استفاده از انتگرال فازی چقدر می‌تواند نسبت به روش‌های سنتی و معمولی چون میانگین وزنی موثر واقع شود. البته در شکل ۵-۸ باید توجه کرد که روش SEP در همان ابتدای کار خود به شدت میانگین حرکت عامل‌ها را کاهش داده ولی به دلیل ماهیت الگوریتم SEP اشباع جداول الگوریتم توانایی ادامه‌ی سرشکن کردن بیشتر میانگین حرکت عامل‌ها را ندارد. میانگین نتایج این قسمت را می‌توان در جدول ۵-۳ خلاصه کرد.

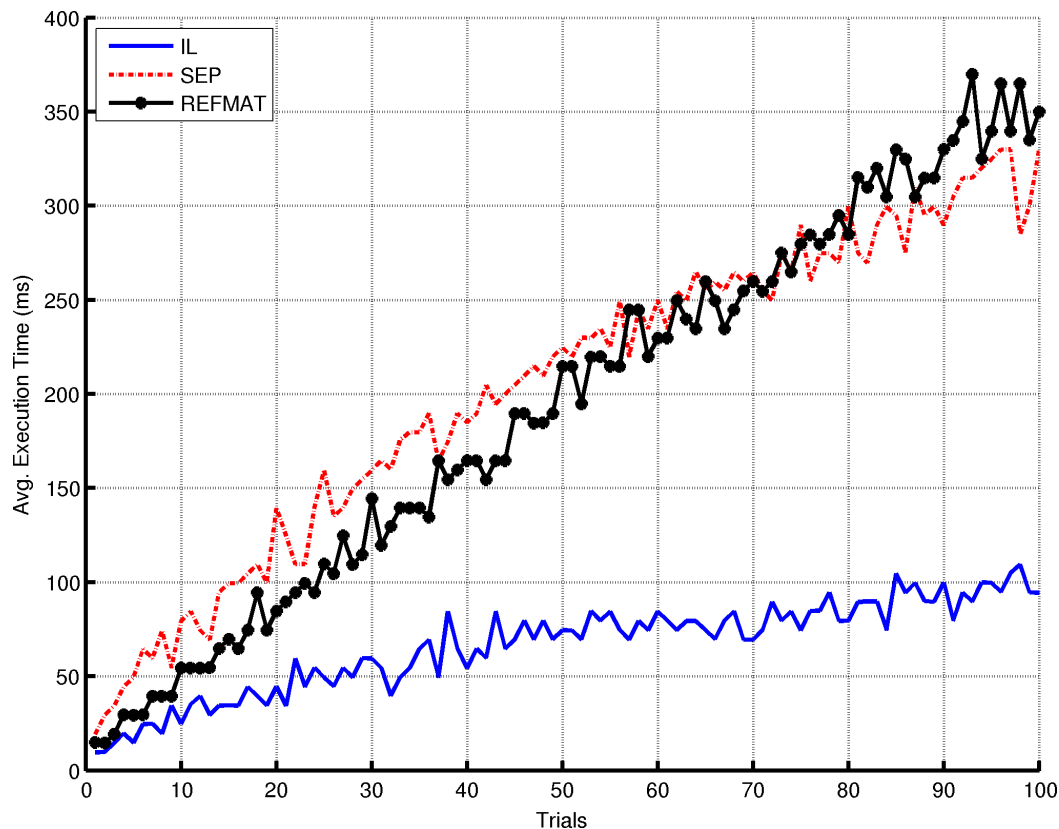
مقایسه در سرعت اجرا: در شکل ۵-۹ میانگین زمانی ۲۰ اجرای مستقل برحسب میلی‌ثانیه به ازای هریک از تعداد تلاش‌ها آورده شده است. همان‌طور که در این شکل مشاهده می‌شود الگوریتم IL دارای حداکثر سرعت



شکل ۵-۸: مقایسه در سرعت و کیفیت یادگیری با تابع  $\varepsilon$ -حریصانه در محیط پلکان مارپیچ

جدول ۵-۳: مقایسه در میزان درصد بهبود کیفیت یادگیری در محیط پلکان مارپیچ با تابع  $\varepsilon$ -حریصانه

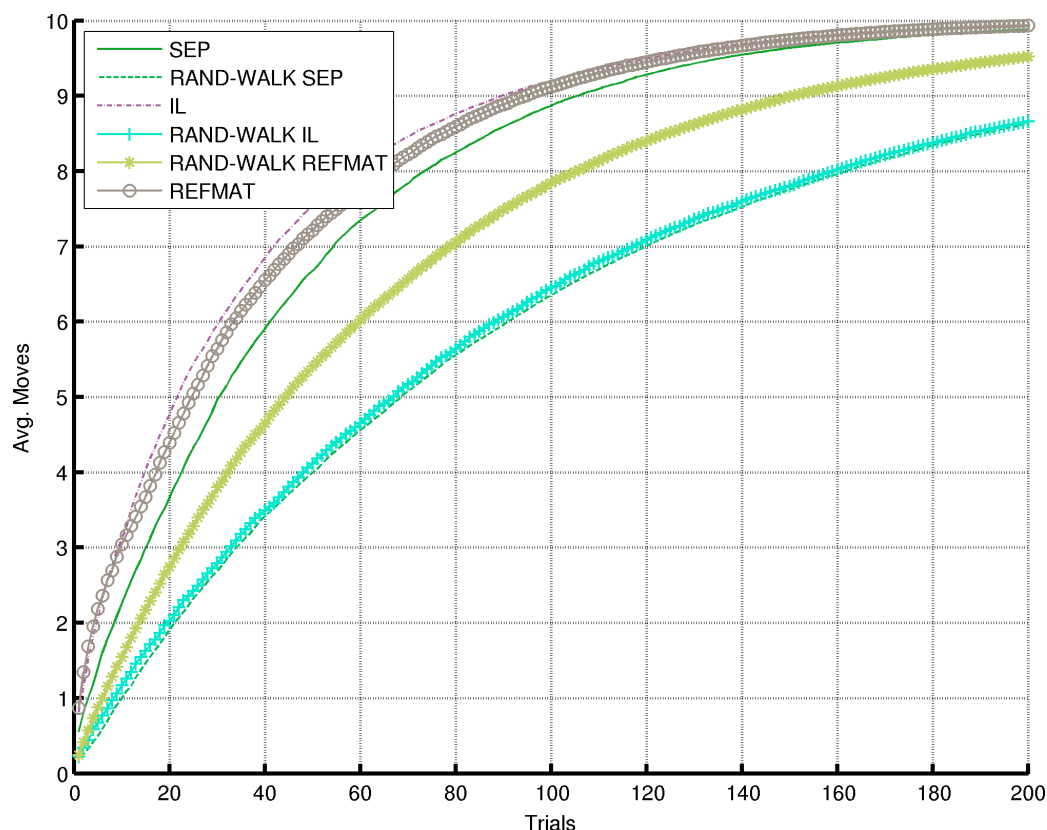
		REFMAT					
	IL	SEP	wsum	fci-mean	fci-max	fci-k-mean	fci-const-one
IL	0.0						
SEP	-3.0	0.0					
wsum	-7.3	-4.4	0.0				
fci-mean	0.6	3.7	8.5	0.0			
fci-max	12.2	15.6	20.9	11.5	0.0		
fci-k-mean	24.0	27.8	33.7	23.2	10.5	0.0	
fci-const-one	34.5	38.7	45.1	33.7	20.0	8.5	0.0



شکل ۵-۹: مقایسه در سرعت اجرای روش‌ها به ازای تعداد تلاش‌های متفاوت برحسب میلی‌ثانیه با تابع  $\varepsilon$ -حریصانه در محیط پلکان مارپیچ

اجرا می‌باشد زیرا که هیچ سربار محاسباتی یادگیری مشترک را ندارد؛ هدف یادگیری اشتراکی این است که می‌خواهد در ازای یک سری سربار محاسباتی کیفیت و سرعت «یادگیری» عامل‌ها را افزایش دهد. با در نظر داشتن این موضوع همانطور که قبلاً دیدیم روش پیشنهادی سرعت و کیفیت یادگیری را بیشتر از روش SEP افزایش می‌دهد و در اینجا نیز می‌بینیم که دارای سرعت اجرای بیشتری نسبت به روش SEP می‌باشد که نشان از بهینه‌گی روش پیشنهادی نسبت به روش SEP می‌دهد.

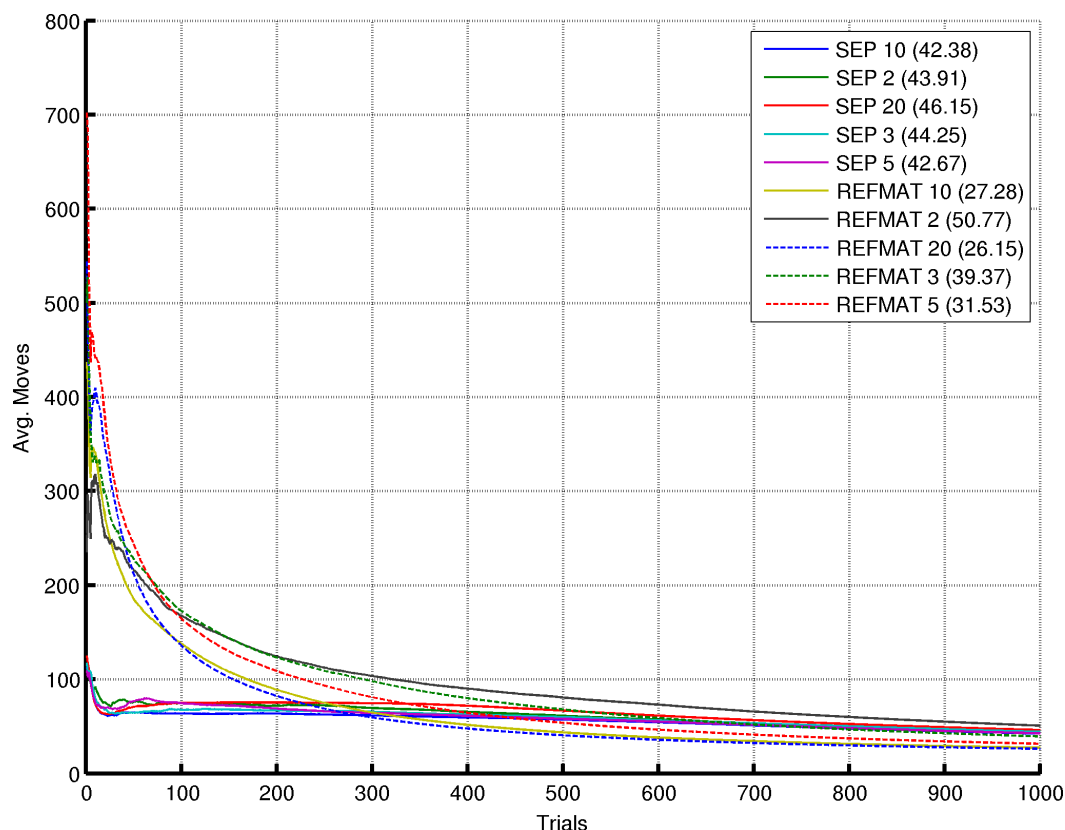
**مقایسه در میزان باروری:** در شکل ۵-۱۰ میزان باروری IL از کلیه روش‌ها بهتر بوده (با اندک اختلاف نسبت روش پیشنهادی) ولی همچنان باروری روش پیشنهادی از روش SEP بیشتر بوده است و همچون آزمایش مشابه با تابع بولتزمن در اینجا نیز نشان داده شده است که روش SEP کاملاً وابسته است به این‌که در هنگام انتخاب عمل بر اساس دانش عامل عمل شود و اگر عامل بدون در نظر گرفتن دانش عامل حرکتی اتخاذ کند میزان باروری عامل بشدت تحت تاثیر قرار می‌گیرد در حالی که در روش پیشنهادی در شرایط یکسان از کلیه الگوریتم‌ها میزان باروری بیشتری دارد.



شکل ۵-۱۰: نمودار باروری الگوریتم‌ها مختلف با تابع  $\varepsilon$ -حریصانه در محیط پلکان مارپیچ

مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری: همان‌طور که در شکل ۵-۱۱ آمده است، روش پیشنهادی و روش SEP به ازای تعداد عامل‌های ۲، ۳، ۵، ۱۰ و ۲۰ عدد به تعداد ۲۰ بار اجرا درآمده و میانگین اجراها به نمودار کشیده شده است. همان‌طور که می‌بینیم روش SEP در زمانی ۲۰ عامل در حال یادگیری و اشتراک گذاری دانش‌های خود هستند نسبت به زمانی که فقط ۲ عامل در حال تعامل مشارکتی با محیط هستند ۸-٪ در خروجی الگوریتم تاثیر منفی داشته است؛ بدین معنی که در زمانی که از تابع  $\varepsilon$ -حریصانه استفاده شود روش SEP به افزایش تعداد عامل فقط منجر به بدتر شدن عملکرد عامل‌ها در یادگیری مشارکتی می‌شود. این در حالی است که در همین شرایط میزان بهبود نتیجه‌ی روش پیشنهادی ۹۲٪ می‌باشد که نشان می‌دهد روش پیشنهادی در ازای افزایش تعداد عامل‌ها به دلیل اینکه دانش جمعی نیز افزایش می‌یابد کیفیت خروجی آن نیز بطور چشم‌گیری بهتر می‌شود. در حالی که در روش SEP اگر کار نتایج بدتر نشود بهتر نمی‌شود که از ضعف بزرگ روش SEP خبر می‌دهد.

نتیجه‌گیری: نتیجه‌ای که از مقایسه‌ی روش پیشنهادی در هر چهار مقایسه‌ی بالا می‌توان گرفت همچون نتیجه‌ای که از نتایج تابع بولتزمن، روش پیشنهادی بهبود چشم‌گیری به روش SEP در محیط پلکان مارپیچ و سیاست انتخاب عمل  $\varepsilon$ -حریصانه داده است.



شکل ۵-۱۱: مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری با تابع  $\epsilon$ -حریصانه در محیط پلکان مارپیچ

جدول ۵-۴: مقایسه در سرعت و کیفیت یادگیری نسبت کیفیت نتیجه‌ی حاصل از تابع  $\epsilon$ -حریصانه نسبت به تابع بولتزمن

		Boltzmann	
		SEP	REFMAT
$\epsilon$ -greedy	SEP	7.27	9.42
	REFMAT	5.20	6.79

مقایسه‌ی بین نتایج حاصل از سیاست انتخاب عمل بولتزمن و  $\epsilon$ -حریصانه

در حالت کلی در محیط پلکان مارپیچ تابع بولتزمن نتایج یکنواثر و پایدارتری<sup>۱</sup> نسبت به تابع  $\epsilon$ -حریصانه از خود نشان داد و در هر دوی این توابع روش پیشنهادی نتیجه‌ی بهتری نسبت به روش SEP ارائه داد. در این قسمت به مقایسه‌ی نتایج بدست آمده توسط هر دو روش در هر دو سیاست انتخاب عمل می‌پردازیم.

مقایسه در سرعت و کیفیت یادگیری: مقایسه‌ی این قسمت را بطور خلاصه می‌توان در جدول ۵-۴ مشاهده کرد. که نسبت کیفیت نتیجه‌ی حاصل از تابع  $\epsilon$ -حریصانه نسبت به تابع بولتزمن همگی بزرگتر از ۱ می‌باشد، که نشان می‌دهد که استفاده از تابع  $\epsilon$ -حریصانه در کیفیت خروجی تأثیری منفی دارد.

<sup>1</sup> Stable

جدول ۵-۵: مقایسه در نسبت میانگین سرعت اجرای حاصل از استفاده تابع  $\varepsilon$ -حریصانه نسبت به تابع بولتزمن

		Boltzmann		
		SEP	REFMAT	IL
$\varepsilon$ -greedy	SEP	1.64	2.05	10.23
	REFMAT	1.72	2.15	10.73
	IL	0.56	0.70	3.49

جدول ۵-۶: مقایسه در نسبت میزان باروری حاصل از استفاده تابع  $\varepsilon$ -حریصانه نسبت به تابع بولتزمن

		Boltzmann		
		SEP	REFMAT	IL
$\varepsilon$ -greedy	SEP	1.08	1.25	1.23
	REFMAT	1.03	1.20	1.18
	IL	1.09	1.27	1.25

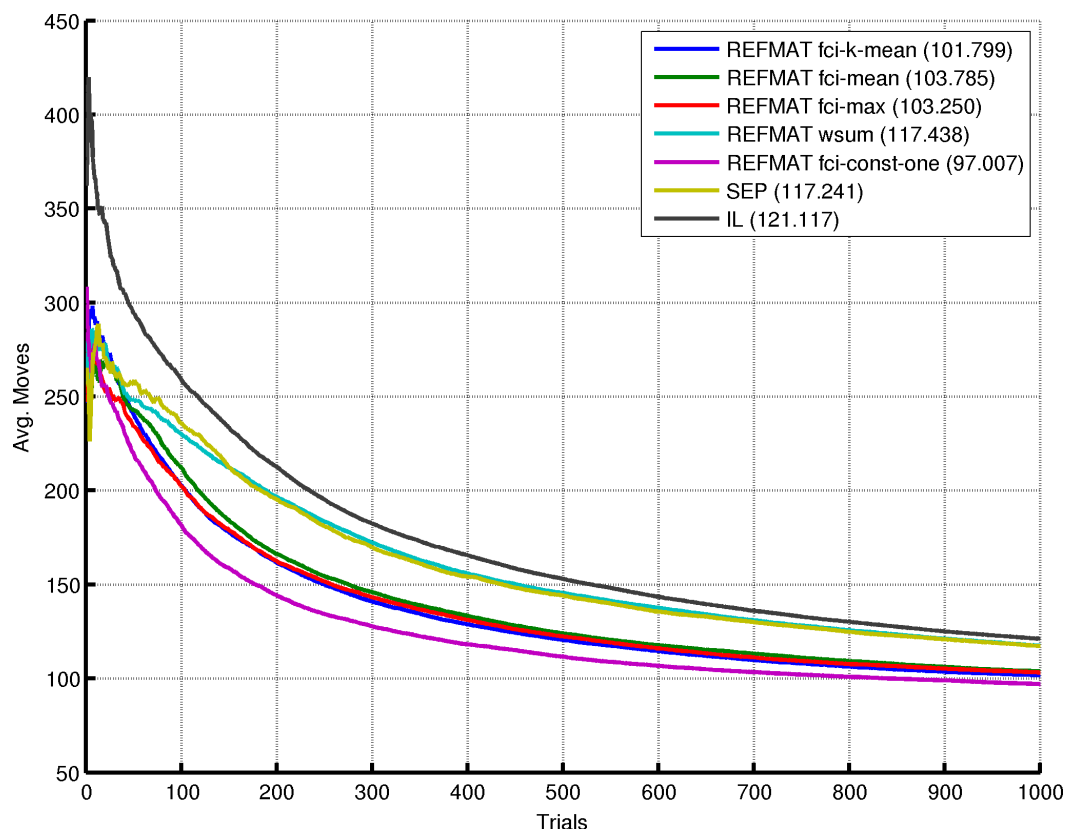
**مقایسه در سرعت اجرا:** در جدول ۵-۵ نسبت میانگین سرعت اجرای روش‌ها آمده است، قطر اصلی این جدول همگی مقادیر بزرگتر از ۱ دارد که نشان می‌دهد هر روش در زمانی که از تابع  $\varepsilon$ -حریصانه استفاده می‌کند زمان بیشتری را تلف می‌کند (صرف جستجوی بی‌مورد محیط می‌کند) نسبت به زمانی که از تابع بولتزمن استفاده می‌کند. این مساله نشان می‌دهد که تابع بولتزمن سریع‌تر عامل را به سمت اهداف هدایت می‌کند - که این نکته در قسمت «مقایسه‌ی سرعت و کیفیت یادگیری» نیز قابل استنتاج است.

**مقایسه در میزان باروری:** همانطور که در جدول ۵-۶ آمده است همه‌ی مقادیر نسبت‌ها بیشتر از ۱ می‌باشد که بدین معنی است که استفاده از تابع  $\varepsilon$ -حریصانه با این حال که کیفیت و سرعت یادگیری کمتری نسبت به تابع بولتزمن دارد و عامل‌ها در حالت کلی زمان زیادی صرف گشت و گذار در محیط می‌کند به نسبت باعث باروری بیشتر جدول  $Q$  می‌شود.

**مقایسه تاثیر تعداد عامل‌ها بر میزان کیفیت و سرعت یادگیری:** در جدول ۵-۷ نسبت شیب تاثیر تعداد عامل‌ها بر میزان کیفیت نتیجه‌ی حاصل از تابع  $\varepsilon$ -حریصانه نسبت به تابع بولتزمن آمده است؛ همانطور که مشاهده می‌شود در زمانی که از تابع  $\varepsilon$ -حریصانه استفاده می‌شود در روش پیشنهادی تاثیر تعداد عامل‌ها به مراتب بیشتر از زمانی است که از تابع بولتزمن استفاده می‌کنیم. این در حالی می‌باشد که در روش SEP اضافه کردن عامل‌ها به محیط تفاوت زیادی در دانش خروجی الگوریتم در هر دو تابع ایجاد نمی‌کند.

جدول ۵-۷: مقایسه نسبت شیب تاثیر تعداد عامل‌ها میزان کیفیت نتیجه‌ی حاصل از تابع  $\varepsilon$ -حریصانه نسبت به تابع بولتزمن

		Boltzmann	
		SEP	REFMAT
$\varepsilon$ -greedy	SEP	0.59	0.09
	REFMAT	73.02	10.67



شکل ۵-۱۲: مقایسه در سرعت و کیفیت یادگیری در محیط صید و صیاد با تابع بولتزمن در محیط صید و صیاد

### ۵-۳-۲ مقایسه در محیط صید و صیاد

همانند مقایسه در محیط پلکان مارپیچ، آزمایش‌های مربوط به این قسمت در ۴ بخش صورت گرفته است؛ ۱. مقایسه در سرعت و کیفیت یادگیری، ۲. مقایسه در پیچیدگی زمانی، ۳. مقایسه در میزان باروری، ۴. مقایسه تاثیر تعداد عامل‌ها بر میزان کیفیت و سرعت یادگیری.

#### سیاست انتخاب عمل «بولتزمن»

مقایسه در سرعت و کیفیت یادگیری: نتایج حاصل از اجرای الگوریتم‌ها در محیط صید و صیاد در شکل ۵-۱۲ آمده است. همانطور که مشاهده می‌شود روش SEP دارای ۳٪ بهبود نسبت به IL می‌باشد ولی همانطور که مشاهده می‌کنیم روش پیشنهادی بهبود چشمگیری نسبت به روش SEP دارد و در بدترین حالت در صورت



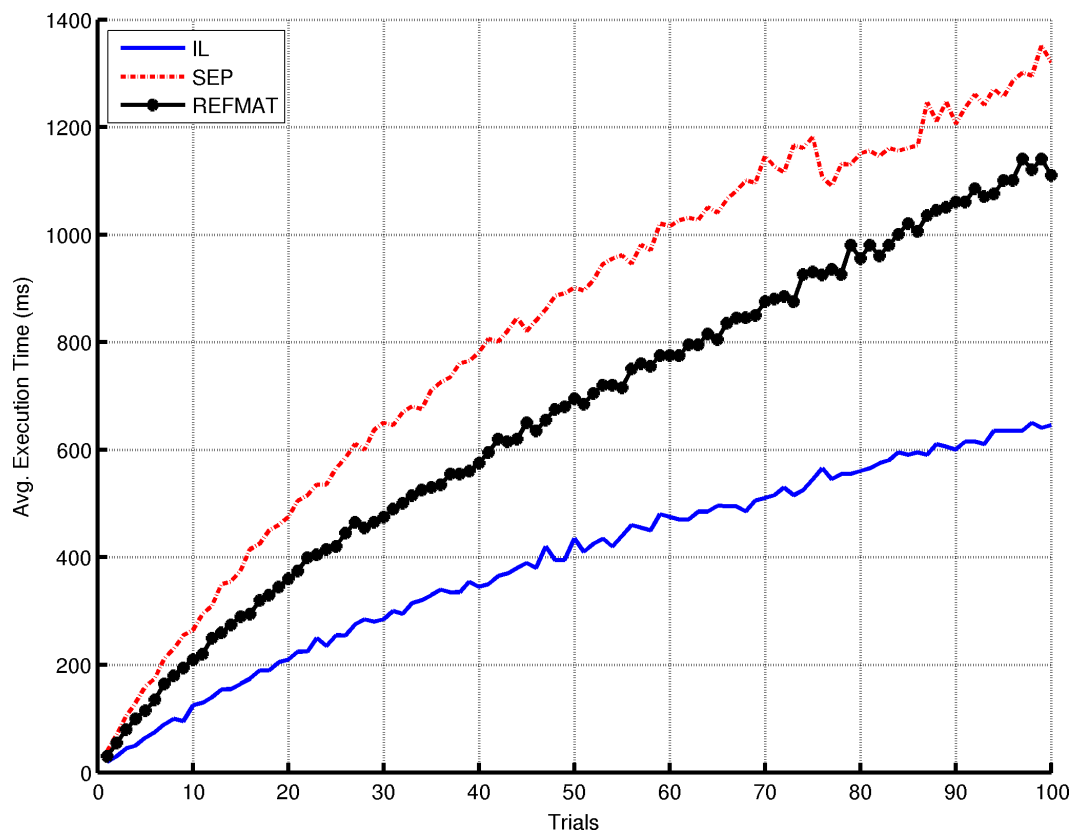
جدول ۵-۸: مقایسه در میزان درصد بهبود کیفیت یادگیری در محیط صید و صیاد با تابع بولترمن

			REFMAT				
	IL	SEP	wsum	fci-mean	fci-max	fci-k-mean	fci-const-one
IL	0.0						
SEP	3.3	0.0					
wsum	3.1	-0.2	0.0				
fci-mean	16.7	13.0	13.2	0.0			
fci-max	17.3	13.5	13.7	0.5	0.0		
fci-k-mean	19.0	15.2	15.4	2.0	1.4	0.0	
fci-const-one	24.9	20.9	21.1	7.0	6.4	4.9	0.0

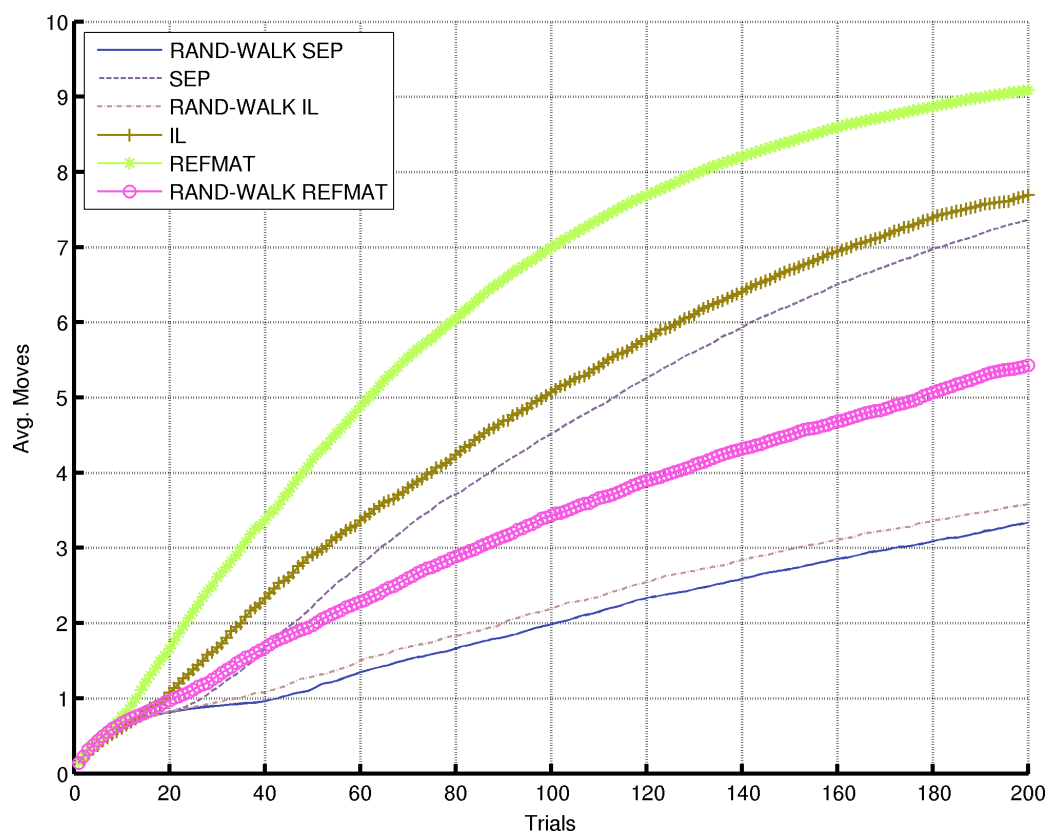
استفاده از تابع Mean به عنوان مدل‌کننده‌ی تابع  $g(\cdot)$  حدود ۱۷٪ نتایج بهبود می‌یابد و در صورتی‌که بخواهیم با استفاده از تابع Const-One به صورت حریصانه عمل کنیم (یعنی در هنگام ترکیب دانش عامل‌ها، فقط دانش عاملی را در نظر بگیریم که با توجه به معیار خبرگی دارای بیشترین خبرگی می‌باشد) مطابق نتایج بدست آمده در محیط قبلی بهترین نتیجه‌ی ممکن یعنی حدود ۲۵٪ بهبود را بدست می‌آوریم؛ نتایج این قسمت را می‌توان در جدول ۵-۸ خلاصه کرد.

**مقایسه در سرعت اجرا:** در شکل ۵-۱۳ همانند شرایط طرح شده در محیط پلکان مارپیچ، سرعت اجرای الگوریتم‌ها آورده شده است. نتایج بدست آمده در محیط صید و صیاد همانند محیط پلکان مارپیچ نشان می‌دهد که روش پیشنهادی نسبت به روش SEP از سرعت اجرای بیشتری برخوردار است.

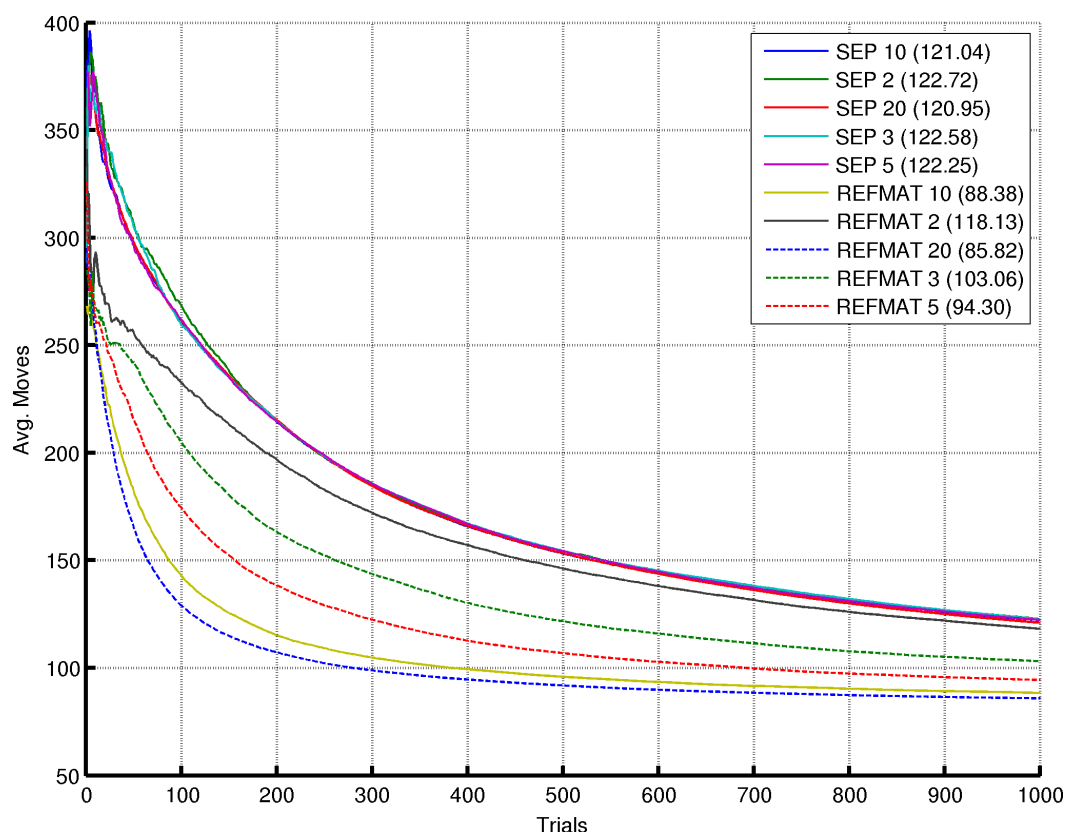
**مقایسه در میزان باروری:** همانطور که در شکل ۵-۱۴ مشاهده می‌کنیم روش معرفی شده در زمانی که به صورت تصادفی اقدام به انتخاب عمل می‌کند بیشتر از زمانی که IL و SEP با بصورت تصادفی اقدام به انتخاب عمل می‌کند جدول Q را بارور می‌کند که از قدرت روش ارائه شده خبر می‌دهد. همچنین در مورد روش SEP می‌بینیم که در زمانی که بصورت تصادفی اقدام به عمل می‌کند باروری کمتری نسبت به روش پیشنهادی و IL دارد؛ یعنی میزان باروری روش SEP وابستگی زیادی به سیاست انتخاب عمل دارد و در صورت نداشتن سیاست انتخاب عمل خاصی شدت عملکردش کاسته می‌شود ولی در روش پیشنهادی میزان این وابستگی از شدت کمتری برخوردار است که از دیگر امتیازات مثبت روش پیشنهادی می‌باشد - همانند نتایج حاصله در محیط پلکان مارپیچ.



شکل ۵-۱۳: مقایسه در سرعت اجرای روش‌ها به ازای تعداد تلاش‌های متفاوت برحسب میلی‌ثانیه با تابع بولتزمن در محیط صید و صیاد



شکل ۵-۱۴: نمودار باروری الگوریتم‌ها مختلف با تابع بولتزمن در محیط صید و صیاد



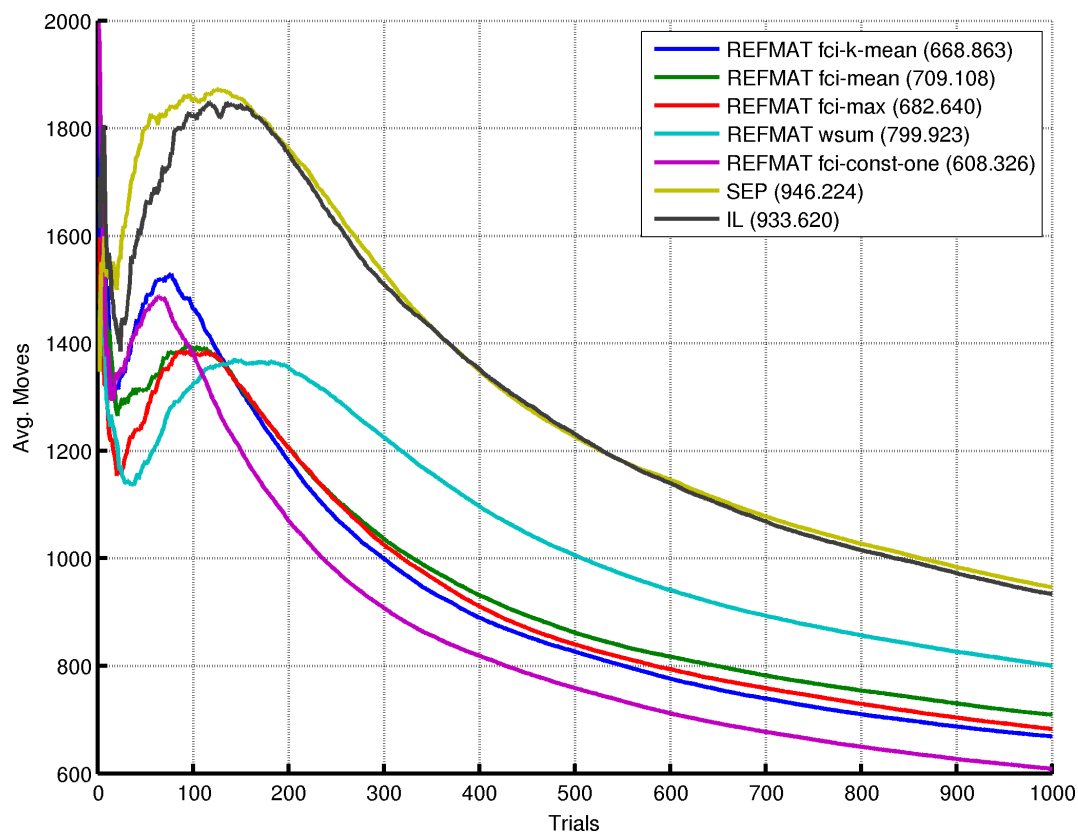
شکل ۵-۱۵: مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری با تابع بولتزمن در محیط صید و صیاد

مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری: همان‌طور که در شکل ۵-۱۵ آمده است، روش SEP در زمانی ۲۰ عامل در حال یادگیری و اشتراک گذاری دانش‌های خود هستند نسبت به زمانی که فقط ۲ عامل در حال تعامل مشارکتی با محیط هستند فقط ۲٪ در خروجی الگوریتم تاثیر مثبت داشته است. این در حالی است که در همین شرایط میزان بهبود نتیجه‌ی روش پیشنهادی ۳۸٪ می‌باشد. که نشان می‌دهد روش SEP نسبت به افزایش تعداد عامل‌ها رفتار تقریباً خنثی از خود نشان می‌دهد در حالی که روش پیشنهادی در ازای افزایش تعداد عامل‌ها به دلیل اینکه دانش جمعی نیز افزایش می‌یابد کیفیت خروجی آن نیز بهتر می‌شود - دلایل و شهود این مساله همانند شهود مطرح شده در محیط صید و صیاد می‌باشد.

نتیجه‌گیری: نتیجه‌ای که از مقایسه‌ی روش پیشنهادی در هر چهار مقایسه‌ی بالا می‌توان گرفت این است که روش پیشنهادی بهبود چشم‌گیری به روش SEP در محیط صید و صیاد و سیاست انتخاب عمل بولتزمن داده است.

#### سیاست انتخاب عمل « $\varepsilon$ -حریصانه»

مقایسه در سرعت و کیفیت یادگیری: نتایج حاصل از اجرای الگوریتم‌ها در محیط صید و صیاد در شکل ۵-۱۶ آمده است. شرایط این آزمایش به مشابه شرایط آزمایش با تابع بولتزمن می‌باشد.

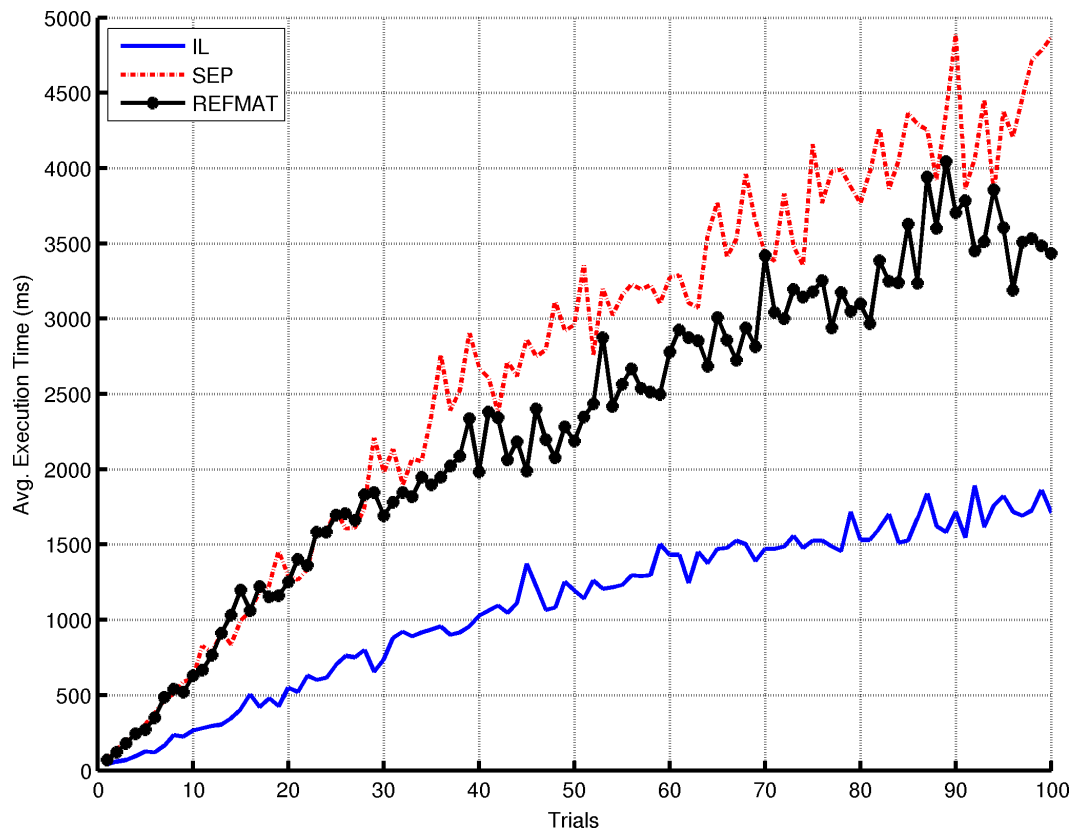


شکل ۵-۱۶: مقایسه در سرعت و کیفیت یادگیری با تابع  $\epsilon$ -حریصانه در محیط صید و صیاد

جدول ۵-۹: مقایسه در میزان درصد بهبود کیفیت یادگیری در محیط صید و صیاد با تابع  $\epsilon$ -حریصانه

	REFMAT						
	IL	SEP	wsum	fci-mean	fci-max	fci-k-mean	fci-const-one
IL	0.0						
SEP	-1.3	0.0					
wsum	16.7	18.3	0.0				
fci-mean	31.7	33.4	12.8	0.0			
fci-max	36.8	38.6	17.2	3.9	0.0		
fci-k-mean	39.6	41.5	19.6	6.0	2.1	0.0	
fci-const-one	53.5	55.5	31.5	16.6	12.2	10.0	0.0

همانطور که خلاصه‌ی این نتایج را در جدول ۵-۹ مشاهده می‌کنیم می‌بینیم که همانند نتایج بدست آمده در آزمایش‌های قبلی روش پیشنهادی با حدود ۵۳٪ بهبود نسبت به IL داشته است درحالی که روش SEP با تابع  $\epsilon$ -حریصانه نه تنها نتایج بهبود داده نشده است بلکه حدود ۱٪- بدتر شده است!

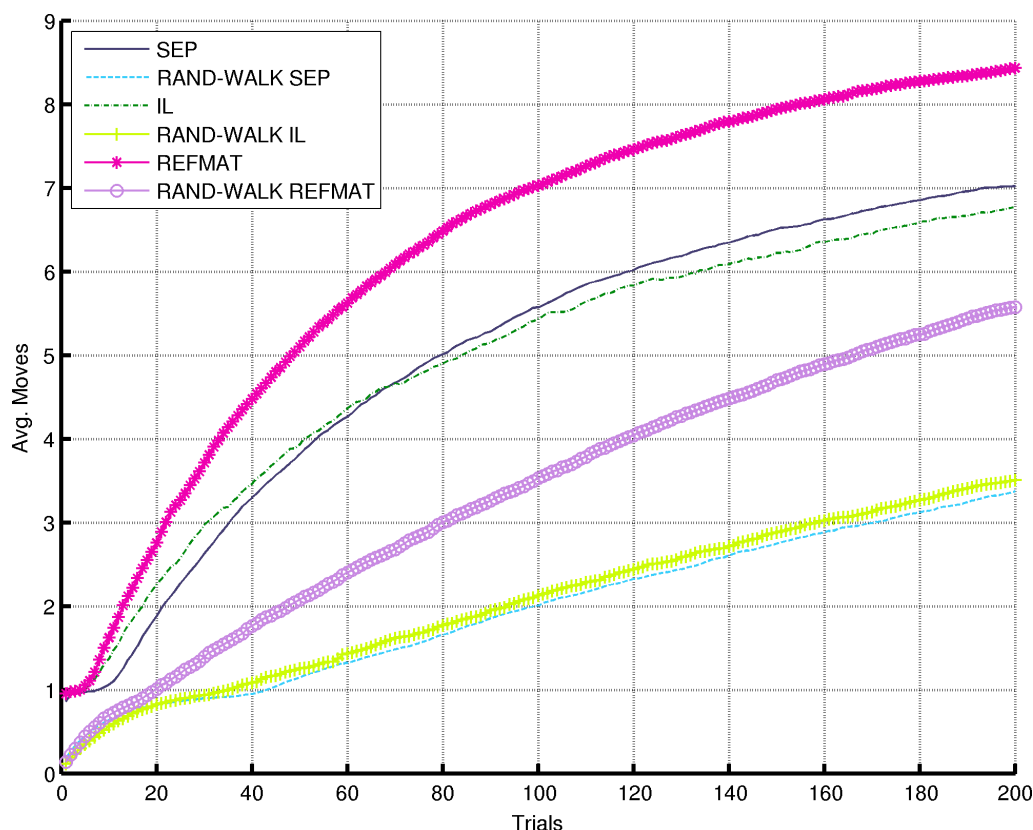


شکل ۵-۱۷: مقایسه در سرعت اجرای روش‌ها به ازای تعداد تلاش‌های متفاوت برحسب میلی‌ثانیه با تابع  $\varepsilon$ -حریصانه در محیط صید و صیاد

مقایسه در سرعت اجرا: در شکل ۵-۱۷ نیز می‌بینیم که در محیط صید و صیاد نیز روش پیشنهادی دارای سرعت اجرای بیشتری نسبت به روش SEP می‌باشد که نشان از بهینه‌گی روش پیشنهادی نسبت به روش SEP می‌دهد.

مقایسه در میزان باروری: در شکل ۵-۱۸ میزان باروری روش پیشنهادی از دیگر روش‌ها بیشتر بوده و همانند آزمایش‌های قبلی در اینجا نیز نشان داده شده است که روش SEP در زمانی که به صورت تصادفی محیط را کاوش کند کمترین باروری را دارد که مطابق دلایل ذکر شده در مقایسه‌ی میزان باروری در آزمایش‌های گذشته این مساله نشان از ضعف بزرگ روش SEP می‌دهد.

مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری: همان‌طور که در شکل ۵-۱۹ آمده است، روش SEP در زمانی ۲۰ عامل در حال یادگیری و اشتراک گذاری دانش‌های خود هستند نسبت به زمانی که فقط ۲ عامل در حال تعامل مشارکتی با محیط هستند ۹۰٪ در خروجی الگوریتم تاثیر منفی داشته است؛ بدین معنی که در زمانی که از تابع  $\varepsilon$ -حریصانه استفاده شود روش SEP به افزایش تعداد عامل فقط منجر به بدتر شدن عملکرد عامل‌ها در یادگیری مشارکتی می‌شود. این در حالی است که در همین شرایط میزان بهبود نتیجه‌ی روش پیشنهادی ۵۵٪ می‌باشد. نشان می‌دهد روش پیشنهادی در ازای افزایش تعداد عامل‌ها به دلیل اینکه دانش جمعی



شکل ۵-۱۸: نمودار باروری الگوریتم‌ها مختلف با تابع  $\epsilon$ -حریصانه در محیط صید و صیاد

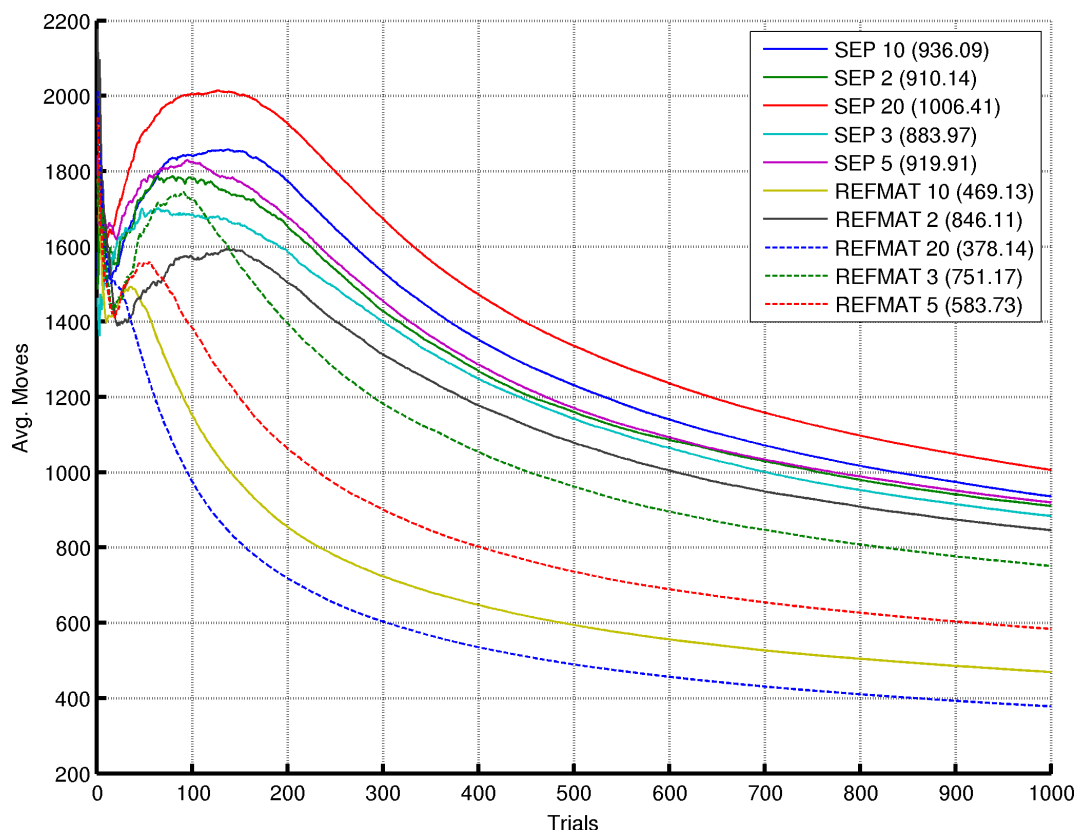
نیز افزایش می‌یابد کیفیت خروجی آن نیز بطور چشم‌گیری بهتر می‌شود. در حالی که در روش SEP اگر کار نتایج بدتر نشود بهتر نمی‌شود که از ضعف بزرگ روش SEP خبر می‌دهد.

**نتیجه‌گیری:** نتیجه‌ای که از مقایسه‌ی روش پیشنهادی در هر چهار مقایسه‌ی بالا می‌توان گرفت همچون نتیجه‌ای که از نتایج تابع بولتزمن، روش پیشنهادی بهبود چشم‌گیری به روش SEP در محیط صید و صیاد و سیاست انتخاب عمل حریصانه داده است.

مقایسه‌ی بین نتایج حاصل از سیاست انتخاب عمل بولتزمن و  $\epsilon$ -حریصانه

در حالت کلی همانند محیط پلکان مارپیچ در محیط صید و صیاد تابع بولتزمن نتایج یکنواتر و پایدارتری نسبت به تابع  $\epsilon$ -حریصانه از خود نشان داد و در هر دوی این توابع روش پیشنهادی نتیجه‌ی بهتری نسبت به روش SEP ارائه داد. در این قسمت به مقایسه‌ی نتایج بدست آمده توسط هر دو روش در هر دو سیاست انتخاب عمل می‌پردازیم.

مقایسه در سرعت و کیفیت یادگیری: مقایسه‌ی این قسمت را بطور خلاصه می‌توان در جدول ۵-۱۰ دید. که نسبت کیفیت نتیجه‌ی حاصل از تابع  $\epsilon$ -حریصانه نسبت به تابع بولتزمن همگی بزرگتر از ۱ می‌باشد، که نشان



شکل ۵-۱۹: مقایسه تاثیر تعداد عامل‌ها میزان کیفیت و سرعت یادگیری با تابع  $\epsilon$ -حریصانه در محیط صید و صیاد

جدول ۵-۱۰: مقایسه در سرعت و کیفیت یادگیری نسبت کیفیت نتیجه‌ی حاصل از تابع  $\epsilon$ -حریصانه نسبت به تابع بولتزمن

		Boltzmann	
		SEP	REFMAT
$\epsilon$ -greedy	SEP	8.07	9.75
	REFMAT	5.19	6.27

می‌دهد که استفاده از تابع  $\epsilon$ -حریصانه در کیفیت خروجی تاثیری منفی دارد.

مقایسه در سرعت اجرا: در جدول ۵-۱۱ نسبت میانگین سرعت اجرای روش‌ها آمده است، که نشان می‌دهد هر روش در زمانی که از تابع  $\epsilon$ -حریصانه استفاده می‌کند زمان بیشتری را تلف می‌کند نسبت به زمانی که از تابع بولتزمن استفاده می‌کند. همانند آنچه که در محیط پلکان مارپیچ مشاهده کردیم در اینجا نیز تابع بولتزمن سریع‌تر عامل را به سمت اهداف هدایت می‌کند.

مقایسه در میزان باروری: همانطور که در جدول ۵-۱۲ آمده است استفاده از تابع  $\epsilon$ -حریصانه به نسبت تابع بولتزمن باعث باروری بیشتر جدول  $Q$  می‌شود.

جدول ۵-۱۱: مقایسه در نسبت میانگین سرعت اجرا حاصل از استفاده تابع  $\varepsilon$ -حریصانه نسبت به تابع بولتزمن

		Boltzmann		
		SEP	REFMAT	IL
$\varepsilon$ -greedy	SEP	3.27	4.10	6.95
	REFMAT	2.74	3.44	5.83
	IL	1.31	1.65	2.79

جدول ۵-۱۲: مقایسه در نسبت میزان باروری حاصل از استفاده تابع  $\varepsilon$ -حریصانه نسبت به تابع بولتزمن

		Boltzmann		
		SEP	REFMAT	IL
$\varepsilon$ -greedy	SEP	1.19	0.82	1.07
	REFMAT	1.49	1.03	1.35
	IL	1.17	0.80	1.05

مقایسه تاثیر تعداد عامل‌ها بر میزان کیفیت و سرعت یادگیری: همانطور که در جدول ۵-۱۳ مشاهده می‌شود در زمانی که از تابع  $\varepsilon$ -حریصانه استفاده می‌شود در روش پیشنهادی تاثیر تعداد عامل‌ها به مراتب بیشتر از زمانی است که از تابع بولتزمن استفاده می‌کنیم. این در حالی می‌باشد که در روش SEP اضافه کردن عامل‌ها به محیط نه تنها به بهبود دانش خروجی الگوریتم کمکی نمی‌کند بلکه نتایج را بدتر نیز می‌کند!

نتیجه‌گیری: با مقایسه‌ی بین تاثیر توابع حریصانه و بولتزمن در خروجی الگوریتم‌ها به این نتیجه می‌توان رسید که تابع بولتزمن رفتاری مطمئن‌تر دارد و باعث می‌شود که روش‌ها سریع‌تر همگرا شود.

#### ۵-۴ بررسی تاثیر تعداد نواحی محیط در کیفیت و سرعت یادگیری عامل‌ها در روش پیشنهادی

همانطور که در تعریف ۴-۱ آورده شده است، بنا به معیار خبرگی معرفی شده در این پژوهش باید محیط به تعدادی ناحیه افزایش شود و سپس میزان حضور عامل در هر ناحیه را سنجیده و خبرگی عامل معکوسی از میزان حضور عامل در این نواحی می‌باشد. لذا ضروری است که در این قسمت به بررسی تاثیر تعداد نواحی محیط در کیفیت و سرعت یادگیری عامل‌ها در روش پیشنهادی بپردازیم.



جدول ۵-۱۳: مقایسه در نسبت شیب تاثیر تعداد عامل‌ها میزان کیفیت نتیجه‌ی حاصل از تابع  $\varepsilon$ -حریصانه نسبت به تابع بولتزمن

		Boltzmann	
		SEP	REFMAT
$\varepsilon$ -greedy	SEP	-2.52	-0.07
	REFMAT	379.32	10.65

#### ۵-۴-۱ محیط پلکان مارپیچ

ما محیط پلکان مارپیچ را به ۶ ناحیه‌ی مختلف با اندازه‌های  $1 \times 1$ ،  $2 \times 2$ ،  $6 \times 6$  (کل محیط) تقسیم‌بندی کرده‌ایم و همان‌طور که در شکل ۵-۲۰ آمده است اندازه‌ی این نواحی در کیفیت و سرعت یادگیری روش پیشنهادی تفاوتی ایجاد نمی‌کند و می‌توان برای کل محیط را یک ناحیه فرض کرد و میزان خبرگی کلی عامل برابر می‌شود با تعداد گام‌هایی که عامل برای رسیدن به هدف طی می‌کند.

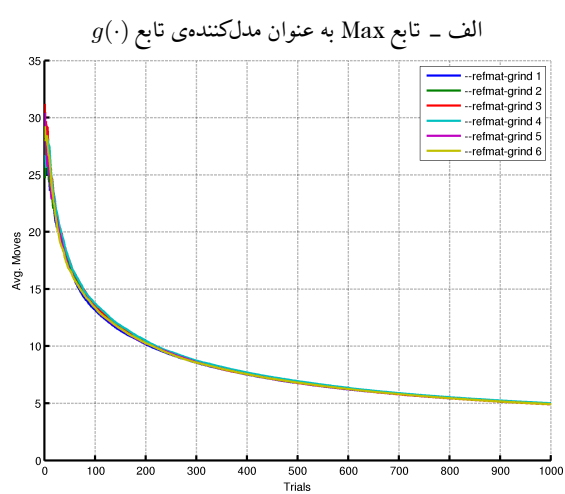
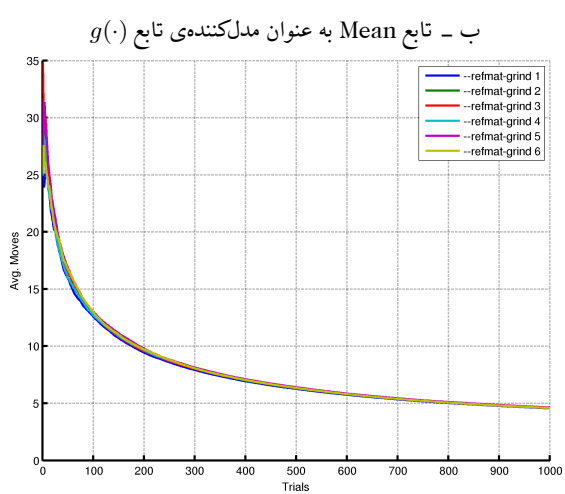
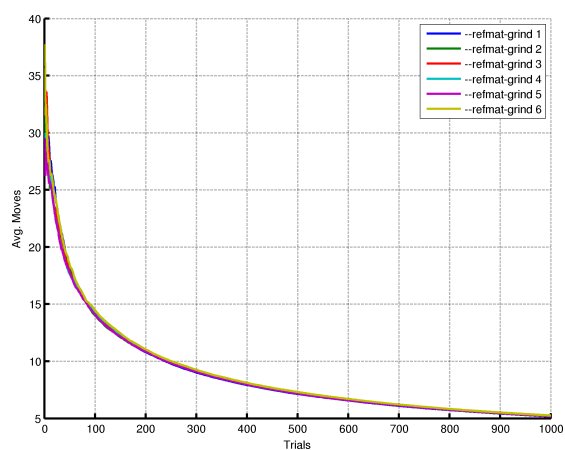
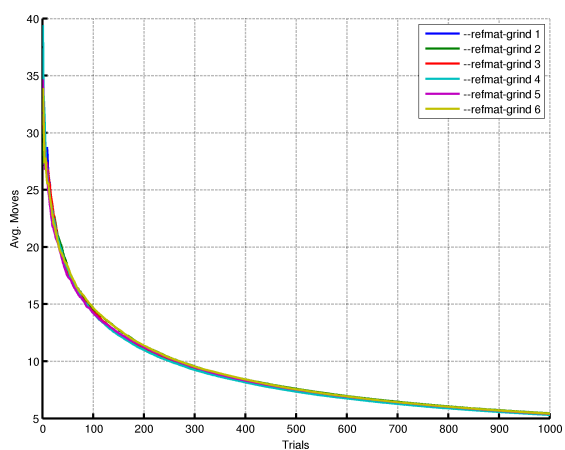
#### ۵-۴-۲ محیط پلکان صید و صیاد

همانند محیط پلکان مارپیچ را به چند ناحیه‌ی مختلف با اندازه‌های  $1 \times 1$ ،  $2 \times 1$ ،  $17 \times 1$  (کل محیط) تقسیم‌بندی کرده‌ایم و همان‌طور که در شکل ۵-۲۱ آمده است همچون محیط پلکان مارپیچ اندازه‌ی این نواحی در کیفیت و سرعت یادگیری روش پیشنهادی تفاوتی ایجاد نمی‌کند.

#### ۵-۵ بررسی تاثیر استفاده از انتگرال فازی در بهبود دانش جمعی

همان‌طور که در بخش‌های قبلی دیدیم، استفاده از توابع  $g(\cdot)$  مختلف در انتگرال فازی نتایج روش پیشنهادی موثر واقع شده‌اند و در نهایت با استفاده از تابع Const-One که باعث می‌شود انتگرال فازی به یک عملگر بیشینه‌گیری تبدیل شود و بطور حریصانه در هر مرحله‌ی به اشتراک‌گذاری دانش، فقط دانش عاملی را که از دیگر عامل‌ها خبره‌تر است به عنوان دانش جمعی خروجی الگوریتم در نظر بگیرد؛ در این بخش می‌خواهیم این موضوع را بررسی کنیم که تاثیر انتگرال فازی در بهبود دانش جمعی چگونه می‌باشد؟ به عنوان مثال آیا همیشه استفاده از انتگرال فازی در بهبود دانش جمعی موثر واقع است؟

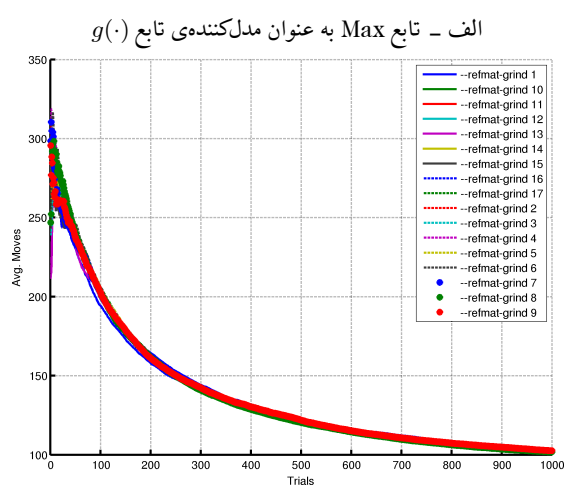
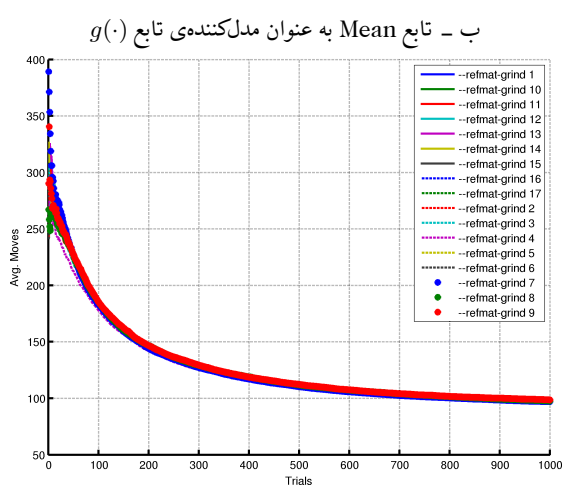
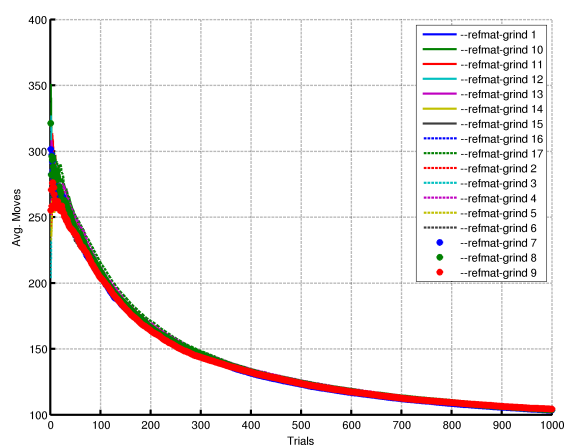
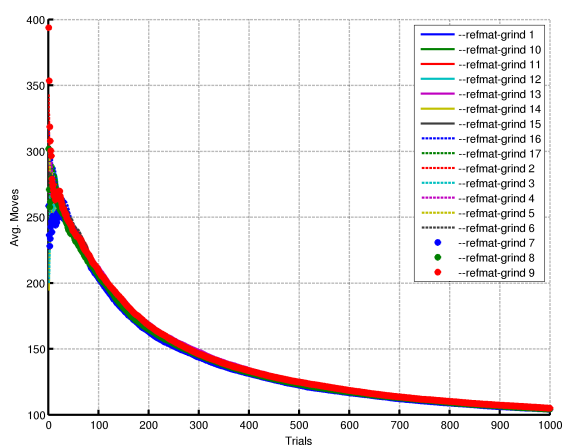
در پاسخ به این سوالات، ما روش SEP را با این بار با استفاده از انتگرال فازی مورد آزمون قرار دادیم که همان‌طور که در شکل ۵-۲۲ مشاهده می‌کنیم وقتی در روش SEP بجای میانگین وزنی (روش پیشنهادی SEP) از انتگرال فازی با توابع معرفی شده در الگوریتم‌های ۴-۵ تا ۴-۷ استفاده کنیم خروجی انتگرال فازی کیفیت و سرعت یادگیری را بهبود نمی‌بخشد - گرچه ممکن است در صورت تعریف توابع  $g(\cdot)$  دیگر نتایج بهتری تولید کند؛ ولی در صورتی که از الگوریتم ۴-۴ استفاده کنیم بهترین نتیجه‌ی ممکن را می‌گیریم.



د - تابع Const-One به عنوان مدل‌کنندهی تابع  $g(\cdot)$

ج - تابع K-Mean به عنوان مدل‌کنندهی تابع  $g(\cdot)$

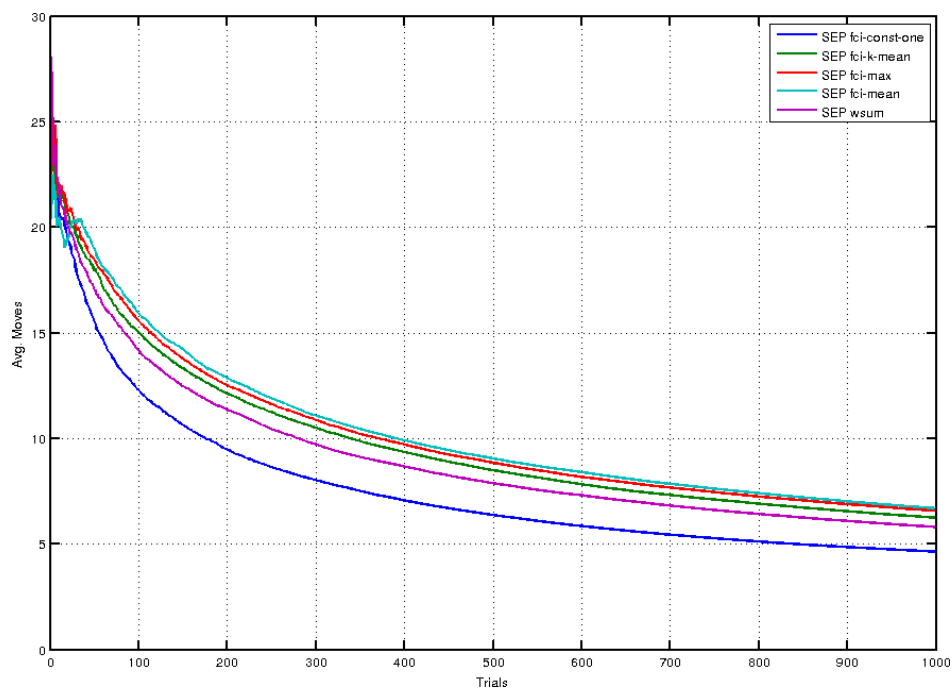
شکل ۵-۲۰: تاثیر ناحیه‌بندی مختلف بروی کیفیت و سرعت یادگیری در محیط پلکان مارپیچ



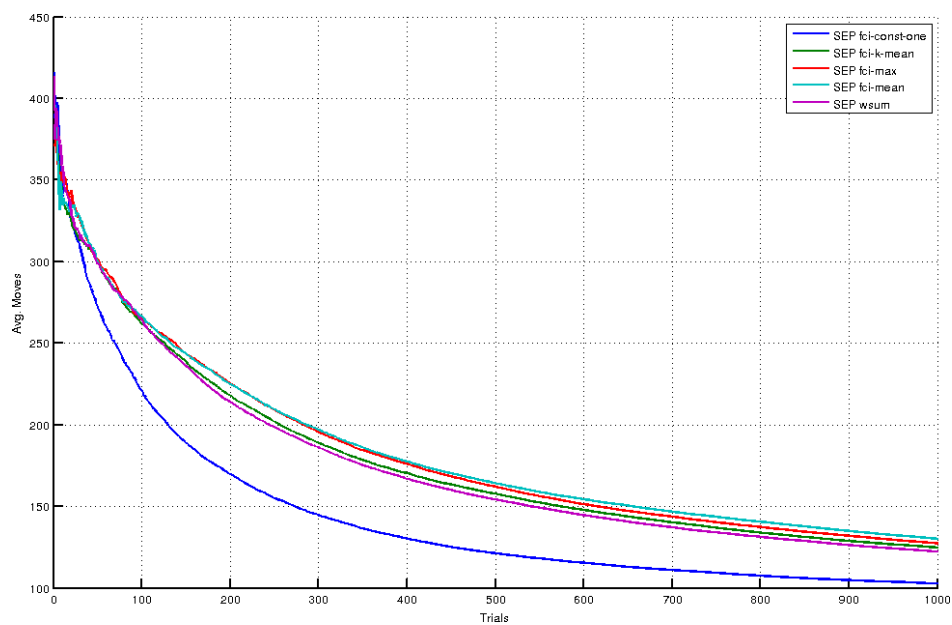
د - تابع Const-One به عنوان مدل‌کنندهی تابع  $g(\cdot)$

ج - تابع K-Mean به عنوان مدل‌کنندهی تابع  $g(\cdot)$

شکل ۵-۲۱: تاثیر ناحیه‌بندی مختلف بروی کیفیت و سرعت یادگیری در محیط صید و صیاد



الف



ب

شکل ۵-۲۲: تاثیر استفاده از انتگرال فازی در روش SEP بروی کیفیت و سرعت یادگیری در محیط پلکان ماریج. الف - استفاده از انتگرال فازی در روش SEP در محیط پلکان ماریج. ب - استفاده از انتگرال فازی در روش SEP در محیط صید و صیاد.

در مورد روش MCE نیز تاثیر استفاده از انتگرال فازی را نیز مورد بررسی قرار دادیم که در شکل ۵-۲۳ نتایج آمده است. همانطور که مشاهده می‌شود انتگرال فازی در روش MCE با توابع  $g(\cdot)$  معرفی شده در این پژوهش (یعنی الگوریتم‌های ۴-۴ تا ۷-۴) باعث بهبود نتایج شده است. این بهبود در محیط پلکان ماریچ به صورت بهبود در سرعت یادگیری و در محیط صید و صیاد در سرعت و کیفیت یادگیری مشاهده می‌شود. همانطور که در بخش‌های قبلی نیز آورده شده است، انتگرال فازی روشی جامع‌تر از میانگین‌گیری وزن‌دار می‌باشد و همانطور که از نتایج آزمایش‌های این فصل می‌توان برداشت کرد با انتخاب و معرفی مناسب تابع  $g(\cdot)$  انتگرال فازی می‌تواند نتایج بهتری نسبت به روش میانگین‌گیری وزن‌دار تولید نماید - این نتیجه‌گیری را می‌توان در بهبود نتایج توسط کلیه توابع  $g(\cdot)$  معرفی شده در این پژوهش بر روی روش پیشنهادی و MCE و همچنین تابع ۴-۴ بروی روش SEP مشاهده کرد، علت بهبود نتایج این روش‌ها با استفاده از انتگرال فازی در قسمت ۴-۴ آورده شده است. همچنین با توجه به نتایج حاصل در شکل ۵-۲۳ باید به این نکته توجه کرد که انتخاب حریصانه‌ی دانش عامل‌ها به عنوان دانش جمعی (یعنی انتخاب توابع Const-One - الگوریتم ۴-۴ به عنوان تابع  $g(\cdot)$ ) همیشه موثر نمی‌باشد.

## ۵-۶ تحلیل نتایج

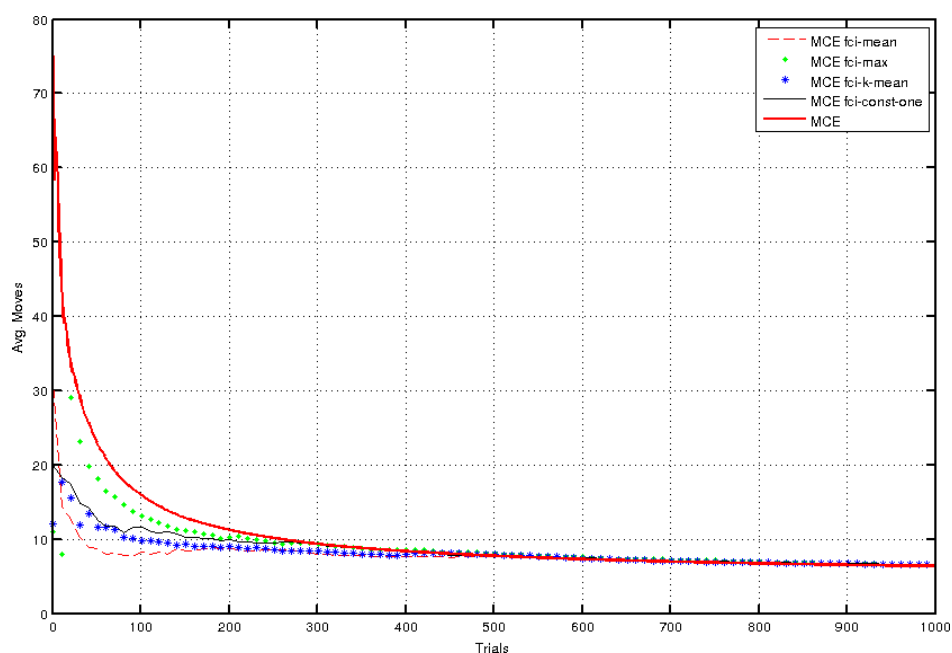
در این قسمت به بررسی و تحلیل نتایج عددی که در بخش‌های قبلی این فصل آورده شده است می‌پردازیم. در تمامی آزمایش‌های انجام شده همانطور که مشاهده می‌شود که روش پیشنهادی از هر لحاظ از روش SEP نتایج بهتری بدست آورده است لذا در این بخش به بررسی علت‌های این برتری می‌پردازیم.

### ۵-۶-۱ مقایسه‌ی روش SEP با روش پیشنهادی

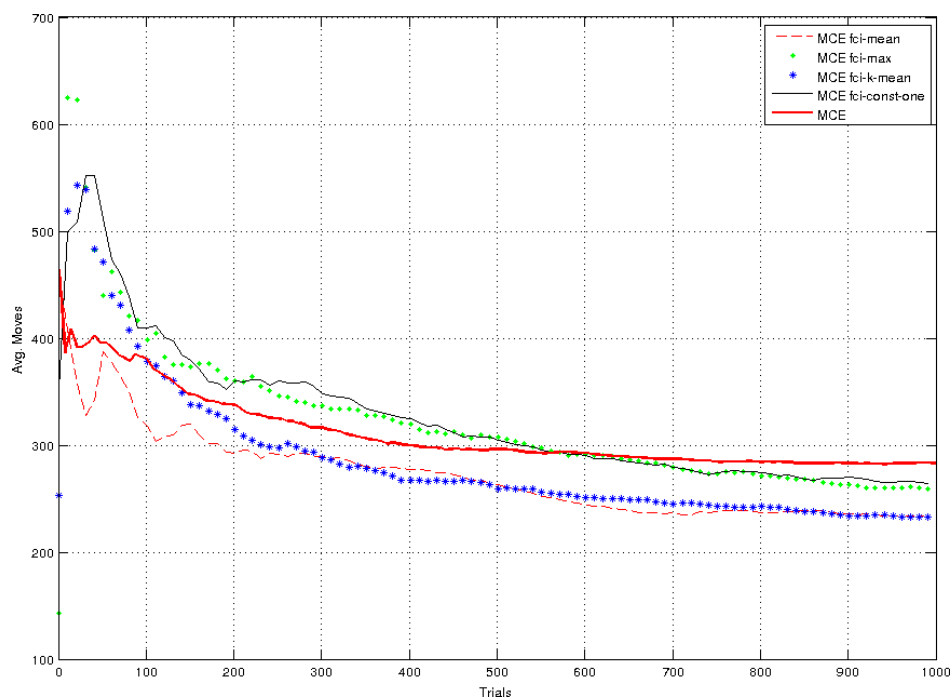
جواب این سوال که «چرا روش پیشنهادی نسبت به روش SEP با اختلاف چشم‌گیری، نتایج بهتری ارائه داده است؟» دلیل زیر می‌باشد:

۱. فقدان وجود پراش<sup>۱</sup> دانشی بین عامل‌ها: برای اینکه ترکیب دانش عامل‌ها موثر واقع شود باید یک پراشی بین دانش عامل‌ها وجود داشته باشد (یعنی عامل‌ها با دانش‌های مختلف از محیط وجود داشته باشد)، اگر همه‌ی عامل‌ها یک دانش را داشته باشند (پراش دانشی کمتری داشته باشند) نباید انتظار داشت که ترکیب دانش عامل‌ها دانشی فراتر از آنچه که در حالت کلی جمع دارند، نتیجه دهد. لذا نمودار شکل ۵-۷ این نتیجه را می‌دهد که روش میرزایی در بین عامل‌ها با تجربیات متفاوت دانش‌های متفاوتی ایجاد نمی‌کند!

<sup>1</sup> Variance



الف



ب

شکل ۵-۲۳: تاثیر استفاده از انتگرال فازی در روش MCE بروی کیفیت و سرعت یادگیری در محیط پلکان مارپیچ. الف - استفاده از انتگرال فازی در روش MCE در محیط پلکان مارپیچ. ب - استفاده از انتگرال فازی در روش MCE در محیط صید و صیاد.

۲. **کیفیت معیارها:** معیارهای معرفی شده توسط میرزایی از کیفیت کافی برخوردار نیستند که نهایتاً در کیفیت نتایج روش ارائه شده توسط ایشان تأثیر می‌گذارند.

#### ۵-۶-۲ مقایسه‌ی تابع بولتزمن و $\varepsilon$ -حریصانه

همانطور که در آزمایش‌ها مشاهده کردیم تابع بولتزمن از هر لحاظ نسبت به تابع  $\varepsilon$ -حریصانه نتایج بهتری تولید می‌کند و شهودی که به نظر می‌رسد دلیل این مساله می‌باشد این است که تابع  $\varepsilon$ -حریصانه در  $\varepsilon\%$  مواقع بصورت کاملاً تصادفی عمل می‌کند. حال فرض کنید عامل بعد از طی چند دوره به دانش نسبتاً خوب نسبت به محیط خود رسیده است، در این صورت در هر حالت بخوبی می‌داند که کدام عمل به صلاح‌تر است، در صورتی که از تابع بولتزمن جهت انتخاب عمل استفاده کند این تابع احتمال انتخاب عملی که در این حالت به صلاح است را بیشتر می‌کند و عملی که تأثیر مخربی بر دستیابی عامل به اهداف خود دارد را احتمال انتخاب کمتری تخصیص می‌دهد؛ به همین دلیل در هر حالت احتمال انتخاب عملی که بهینه است بیشتر است و به مرور که عامل دانش خود را نسبت به محیط کامل‌تر می‌کند این احتمال بیشتر تقویت می‌شود. در حالی که تابع  $\varepsilon$ -حریصانه فارغ از میزان دانش عامل از محیط در هر حالت  $\varepsilon\%$  احتمال دارد که کاملاً تصادفی انتخاب عمل کند و در این انتخاب تصادفی محتمل است که بدترین عمل انتخاب شود و این احتمال‌ها در طور یادگیری ثابت است لذا ممکن است در یک قدمی رسیدن به هدف عامل به صورت تصادفی یک عملی را انتخاب که کند که منجر عدم دستیابی به هدف یا دور شدن از هدف شود، به همین دلیل تابع  $\varepsilon$ -حریصانه نسبت به تابع بولتزمن باعث می‌شود عامل دیرتر به هدف برسد که همین مساله در معیارهای سنجش از قبیل سرعت و کیفیت یادگیری، سرعت اجرا تأثیر منفی می‌گذارد.

#### ۵-۶-۳ بررسی تأثیر تعداد نواحی در کیفیت و سرعت یادگیری در روش پیشنهادی

دلیل این مساله که «با توجه به شکل‌های ۵-۲۰ و ۵-۲۱ تعداد نواحی معیار تعریف شده در تعریف ۴-۱ در سرعت و کیفیت یادگیری در روش پیشنهادی تأثیرگذار نیست.» این است که با توجه به تعریف ۴-۱ یک همبستگی<sup>۱</sup> مثبتی میان میزان ارجاع در نواحی ریز و نواحی درشت وجود دارد، زیرا که هرچه میزان ارجاع به نواحی ریز زیاد باشد همان میزان ارجاع به نواحی درشت (که شامل آن نواحی ریز می‌باشند) نیز زیاد می‌شود و برعکس هرچه میزان ارجاع به نواحی درشت کمتر باشد میزان ارجاع به نواحی زیرمجموعه‌ی آن نواحی درشت نیز کمتر است. بنابراین با توجه نحوه‌ی تعریف معیار مورد استفاده در تعریف ۴-۱ تعداد نواحی در کیفیت و سرعت یادگیری تأثیرگذار نمی‌باشد.

<sup>۱</sup> Correlation

## ۵-۲ نتیجه‌گیری

در فصل به بررسی نتایج آزمایش‌های گوناگونی که برای مقایسه‌ی روش پیشنهادی با روش SEP آورده شده بودند پرداختیم؛ مقایسه‌ها بر اساس ۴ معیار صورت گرفت: ۱. مقایسه در سرعت و کیفیت یادگیری، ۲. مقایسه در پیچیدگی زمانی، ۳. مقایسه در میزان باروری، ۴. مقایسه تاثیر تعداد عامل‌ها بر میزان کیفیت و سرعت یادگیری. همانطور که دیدیم این روش پیشنهادی در هر ۴ معیار نسبت به روش SEP برتری قابل توجهی داشت. همچنین به مقایسه‌ی تاثیر استفاده از سیاست‌های انتخاب عمل بولتزمن و  $\epsilon$ -حریصانه پرداختیم و نشان دادیم که استفاده از بولتزمن نسبت به  $\epsilon$ -حریصانه کیفیت و سرعت یادگیری بسیار بهتری بدست می‌دهد. همچنین نشان دادیم که معیار خبرگی معرفی شده در تعریف ۴-۱ مستقل از تعداد و اندازه‌ی افرازهای محیط می‌باشد، سپس با تحلیل نتایج آزمایش‌ها به مطالب این بخش خاتمه دادیم. در فصل بعدی به جمع‌بندی مطالب و دست‌آوردهای این پژوهش می‌پردازیم.



## فصل ششم

### نتیجه‌گیری و جمع‌بندی

#### ۶-۱ مقدمه

اگر اندکی به مسائلی که افراد انجام می‌دهند و ما آن‌ها را در آن خبره می‌بینیم توجه کنیم، متوجه خواهیم شد که زمانی که فردی در موردی خبره می‌شود بطور طبیعی انرژی نسبتاً کمتری در انجام آن مصرف می‌کند. این معیار همان معیاری است که می‌گوید عاملی در انجام وظیفه‌ای خبره‌تر است که در طی انجام آن انرژی کمتری مصرف کند. این معیار که از فلسفه‌ی بسیار ساده‌ای برخوردار است برخلاف معیارهای گذشته بسیار منعطف می‌باشد زیرا که در تعریف این معیار عبارت «میزان انرژی» می‌تواند تعابیر مختلفی به خود بگیرد و در هر مورد قابل استفاده باشد.

وجود این فلسفه‌انگیزه‌ای شد که در صدد ارائه‌ای معیاری برآیم که نه تنها ساده باشد بلکه در زندگی روزمره ما انسان‌ها هم تجلی داشته باشد. بعد از اندکی تفکر و تفحص در نهایت این معیار چیزی جز معیار «تنبلی» نبود! معیار تنبلی که در این پایان‌نامه با اصطلاح علمی «میزان ارجاع» ارائه شد می‌گوید که «عاملی که برای به نتیجه رساندن فعالیت‌هایش انرژی کمتری صرف کند خبره‌تر است».

در این قسمت به مروری خلاصه بر هرآنچه که در این پژوهش صورت گرفته و ارائه‌ی یک نتیجه‌گیری نهایی حاصل از این پژوهش و همچنین ارائه‌ی مسیر پژوهشی پیشنهادی برای آیندگان این زمینه از یادگیری مشارکتی خواهیم پرداخت.

## ۶-۲ نوآوری‌ها و نتایج کلی پایان‌نامه

در طی این پایان‌نامه معیار جدیدی به نام معیار «میزان ارجاع» ارائه شد که می‌گوید عاملی که کمتر در محیط مورد تعاملش پرسه بزند از خبرگی بیشتری برخوردار است و سپس با استفاده از این معیار خبرگی به سنجش عامل‌های فعال در محیط در هنگام مشارکت در دانش جمعی پرداختیم. در هنگام ترکیب دانش عامل‌ها از انتگرال فازی چوکت استفاده شد.

در طی آزمایش‌ها از میانگین وزنی نیز به جای انتگرال فازی استفاده شد و نشان داده شد که در معیار ارائه شده توسط این پژوهش انتگرال فازی توانایی بهتری نسبت به میانگین وزنی برای بهبود کیفیت و سرعت یادگیری مشارکتی دارد. همچنین از ۴ تابع به عنوان مدل‌کننده‌ی تابع  $g(\cdot)$  استفاده شد، که هرکدام یک دیدگاهی نسبت به نحوه‌ی ترکیب دانش‌های ورودی ارائه می‌دهد. از بین این ۴ تابع، تابع Const-One در کلیه‌ی آزمایش‌ها نسبت به دیگر توابع برتری قابل توجهی از خود نشان داد؛ طبق آنچه که فصول قبلی این پایان‌نامه آورده شده این تابع معادل با حداکثرگیری بروی دانش عامل‌ها بر اساس معیار خبرگی آن‌ها می‌باشد. یعنی اینکه این تابع در واقع در هر ناحیه فقط دانش عاملی را در نظر می‌گیرد از همه خبره‌تر (تنبل‌تر) است که این امر تاییدی بر فرضیه ۴-۱ و متعاقباً تعریف ۴-۱ می‌باشد. همچنین تاثیر استفاده از انتگرال فازی در روش SEP را نیز مورد بررسی قرار دادیم و مشاهده کردیم که انتگرال فازی با توابع  $g(\cdot)$  معرفی شده در این پژوهش در روش SEP تاثیر مثبتی ندارد ولی در این بررسی نیز تابع Const-One بیشترین بهبود را در پی داشت که معادل می‌شود با انتخاب حریصانه‌ی بین دانش عامل‌ها به عنوان دانش جمعی، که یک نقطه‌ی مشترک بین نتایج روش پیشنهادی و این بررسی می‌باشد؛ از طرفی بررسی‌های انجام شده در رابطه با تاثیر انتگرال فازی با توابع  $g(\cdot)$  معرفی شده در این پژوهش در روش MCE نشان می‌دهد که انتگرال فازی موثر واقع شده است، که نشان می‌دهد در صورت انتخاب مناسب تابع  $g(\cdot)$  انتگرال فازی می‌تواند جایگزین بهتری بجای میانگین‌گیری وزنی باشد. همچنین لازم به ذکر است که بررسی‌های انجام شده در استفاده از انتگرال فازی بروی هر سه روش REFMAT و SEP و MCE نشان می‌دهد که انتخاب حریصانه دانش عامل‌ها به عنوان دانش جمعی (انتخاب Const-One به عنوان تابع  $g(\cdot)$ ) همیشه باعث حداکثر شدن سرعت و کیفیت یادگیری نمی‌شود که این نتیجه‌گیری به اهمیت استفاده از انتگرال فازی به عنوان تابع ادغام‌کننده دانش عامل‌ها می‌افزاید.

همچنین در نهایت، در انتهای فصل آزمایش‌ها نشان داده شد که می‌توان معیار خبرگی ارائه شده در تعریف ۴-۱ را به کل محیط خلاصه کرد؛ یعنی عاملی خبره‌تر است که میزان حضور آن در کل محیط کمتر باشد - یعنی با تعداد گام کمتری به اهداف خود برسد. همین نتیجه‌گیری باعث می‌شود که آزمودن دیگر توابع برای مدل کردن  $g(\cdot)$  (مثلاً تابع اندازه‌گیری- $\lambda$  سوگنو) نیازی نباشد.

در این پژوهش تعادلی بین کلی و جزئی نگری به عملکرد عامل‌ها در هنگام ادغام دانش‌های آن‌ها برقرار شد. همچنین تاثیر دیگر روش‌های انتخاب عمل را در ترکیب با معیارهای ارائه شده را مورد بررسی قرار گرفته است و به این نتیجه رسیدیم که تابع بولتزمن نتیجه‌ی با کیفیت‌تری را تولید می‌کند. همچنین دستاوردهای این پژوهش را با در نظر گرفتن ماهیت غیرافزایشی بودن ذات مساله ارائه دادیم.

یکی از مزایای روش پیشنهادی این است که در عین کارایی و قدرت روشی ساده در مفهومی و پیاده‌سازی می‌باشد که این سادگی طبق آنچه که در آزمایش‌ها آمده است نهایتاً منجر شد که روش پیشنهادی از پیچیدگی کمتری برخوردار باشد. از دیگر مزیت روش پیشنهادی کلی بودن فرضیه خبرگی‌ای که این پژوهش بر مبنای آن ارائه شد، می‌باشد که می‌توان آن را به تمامی مسائل یادگیری مشارکتی به راحتی اعمال کرد.

### ۳-۶ راهکارهای آینده و پیشنهادها

همانطور که آزمایش‌ها نشان دادند با توجه به معیار خبرگی ارائه شده در قسمت یادگیری مشارکتی اگر فقط دانش عامل خبره را در نظر بگیریم حداکثر نتیجه‌ی ممکن (در قالب روش پیشنهادی) را خواهیم گرفت. در طی این پژوهش دو مفهوم مهم ارائه شد: ۱. انتگرال فازی چوکت می‌تواند عملکرد بسیار قوی‌ای نسبت به روش‌ها سنتی چون میانگین‌گیری وزنی باشد. ۲. فرضیه خبرگی معرفی شده بخوبی می‌تواند هر نوع معیار خبرگی را توجیه کند. در این پژوهش سعی شده است که حداکثر نتیجه‌ی ممکن حاصل از استفاده از این دو مفهوم باهم را استخراج کنیم ولی پیشنهادها زیر می‌تواند شروع خوبی برای پژوهش‌های آینده در این زمینه باشد.

۱. ارائه‌ی معیار خبرگی جدیدی مبتنی بر فرضیه خبرگی (فرضیه‌ی ۴-۱) معرفی شده در این پژوهش که چهارچوبی کلی جهت تعریف معیارهای خبرگی را تعریف می‌کند؛ سپس آزمایش معیار خبرگی تعریف شده بجهت آزمودن فرضیه خبرگی ارائه شده.

۲. بررسی تاثیر استفاده از انتگرال فازی چوکت در پژوهش‌های گذشته.

۳. بررسی این موضوع که «آیا انتخاب دانش عاملی که از همه خبره‌تر است همیشه باعث حداکثر شدن کیفیت و سرعت یادگیری جمع می‌شود یا خیر؟».

۴. ارائه‌ی روشی جهت وفقی شدن پارامترهای معرفی شده در بخش ۴-۳-۲ با هدف بهبود کیفیت و سرعت یادگیری هرچه بیشتر روش پیشنهادی.

۵. بررسی شرایط و ویژگی‌های توابع  $g(\cdot)$  در حالت کلی.

## مراجع

- [1] E. P. Hajyyar, "Multi-criteria expertness based cooperative learning in multi-agent systems," Master's thesis, Department of Electrical and Computer Engineering, Isfahan University of Technology, Isfahan University of Technology, Isfahan 84156-83111, Iran, 10 2010.
- [2] M. ali mirzaei badizi, "Speed-up cooperative learning in multi-agent systems using shortest experimented path," Master's thesis, Department of Electrical and Computer Engineering, Isfahan University of Technology, Isfahan University of Technology, Isfahan 84156-83111, Iran, 3 2016.
- [3] L. Panait and S. Luke, "Cooperative multi-agent learning: The state of the art," *Autonomous agents and multi-agent systems*, vol. 11, no. 3, pp. 387–434, 2005.
- [4] V. Torra, Y. Narukawa, and M. Sugeno, *Non-Additive Measures*, pp. 3–7. Springer, 2014.
- [5] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *Proceedings of the tenth international conference on machine learning*, pp. 330–337, 1993.
- [6] H. R. Berenji and D. Vengerov, "Cooperation and coordination between fuzzy reinforcement learning agents in continuous state partially observable markov decision processes," in *Fuzzy Systems Conference Proceedings, 1999. FUZZ-IEEE'99. 1999 IEEE International*, vol. 2, pp. 621–627, IEEE, 1999.
- [7] Y. Kuniyoshi, M. Inaba, and H. Inoue, "Learning by watching: Extracting reusable task knowledge from visual observation of human performance," *IEEE transactions on robotics and automation*, vol. 10, no. 6, pp. 799–822, 1994.
- [8] A. Garland and R. Alterman, "Multiagent learning through collective memory," in *Adaptation, Coevolution and Learning in Multiagent Systems: Papers from the 1996 AAAI Spring Symposium*, pp. 33–38, 1996.
- [9] A. Garland and R. Alterman, "Preparation of multi-agent knowledge for reuse," in *Proceedings of the Fall Symposium on Adaptation of Knowledge for Reuse*, vol. 26, p. 33, 1995.
- [10] L. Nunes and E. Oliveira, "On learning by exchanging advice," *arXiv preprint cs/0203010*, 2002.
- [11] L. Nunes and E. Oliveira, "Advice-exchange between evolutionary algorithms and reinforcement learning agents: Experiments in the pursuit domain," in *Adaptive Agents and Multi-Agent Systems II*, pp. 185–204, Springer, 2005.

- [12] M. N. Ahmadabadi, M. Asadpur, S. H. Khodanbakhsh, and E. Nakano, "Expertness measuring in cooperative learning," in *Intelligent Robots and Systems, 2000.(IROS 2000). Proceedings. 2000 IEEE/RSJ International Conference on*, vol. 3, pp. 2261–2267, IEEE, 2000.
- [13] S. M. Eshgh and M. N. AhmadAbadi, "An extension of weighted strategy sharing in cooperative q-learning for specialized agents," in *Neural Information Processing, 2002. ICONIP'02. Proceedings of the 9th International Conference on*, vol. 1, pp. 106–110, IEEE, 2002.
- [14] P. Ritthipravat, T. Maneewarn, J. Wyatt, and D. Laowattana, "Comparison and analysis of expertness measure in knowledge sharing among robots," in *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems*, pp. 60–69, Springer, 2006.
- [15] Y. Yang, Y. Tian, and H. Mei, "Cooperative q learning based on blackboard architecture," in *Computational Intelligence and Security Workshops, 2007. CISW 2007. International Conference on*, pp. 224–227, IEEE, 2007.
- [16] H. Iima and Y. Kuroe, "Reinforcement learning through interaction among multiple agents," in *2006 SICE-ICASE International Joint Conference*, pp. 2457–2462, IEEE, 2006.
- [17] J. Kennedy, "Particle swarm optimization," in *Encyclopedia of machine learning*, pp. 760–766, Springer, 2011.
- [18] M. Yang, Y. Tian, and X. Liu, "Cooperative q-learning based on maturity of the policy," in *2009 International Conference on Mechatronics and Automation*, pp. 1352–1356, IEEE, 2009.
- [19] E. Pakizeh, M. Palhang, and M. M. Pedram, "Multi-criteria expertness based cooperative q-learning," *Applied intelligence*, vol. 39, no. 1, pp. 28–40, 2013.
- [20] D. L. Poole and A. K. Mackworth, *Artificial Intelligence: foundations of computational agents*, ch. 11. Cambridge University Press, 2010.
- [21] S. J. Russell, P. Norvig, J. F. Canny, J. M. Malik, and D. D. Edwards, *Artificial intelligence: a modern approach*, vol. 2. Prentice hall Upper Saddle River, 2003.
- [22] V. Torra and Y. Narukawa, "The interpretation of fuzzy integrals and their application to fuzzy systems," *International journal of approximate reasoning*, vol. 41, no. 1, pp. 43–58, 2006.
- [23] K. Leszczyński, P. Penczek, and W. Grochulski, "Sugeno's fuzzy measure and fuzzy clustering," *Fuzzy Sets and Systems*, vol. 15, no. 2, pp. 147–158, 1985.
- [24] A. F. Tehrani, W. Cheng, and E. Hullermeier, "Preference learning using the choquet integral: The case of multipartite ranking," *IEEE Transactions on Fuzzy Systems*, vol. 20, no. 6, pp. 1102–1113, 2012.
- [25] L. M. De Campos and M. Jorge, "Characterization and comparison of sugeno and choquet integrals," *Fuzzy Sets and Systems*, vol. 52, no. 1, pp. 61–67, 1992.
- [26] M. Grabisch, "Fuzzy integral in multicriteria decision making," *Fuzzy sets and Systems*, vol. 69, no. 3, pp. 279–298, 1995.
- [27] T. Murofushi, M. Sugeno, and M. Machida, "Non-monotonic fuzzy measures and the choquet integral," *Fuzzy sets and Systems*, vol. 64, no. 1, pp. 73–86, 1994.
- [28] M. Grabisch, "The application of fuzzy integrals in multicriteria decision making," *European journal of operational research*, vol. 89, no. 3, pp. 445–456, 1996.
- [29] "Expert - wikipedia." <https://en.wikipedia.org/wiki/Expert>. (Accessed on 11/12/2016).
- [30] E. Schechter, *Handbook of Analysis and its Foundations*, ch. 1, p. 16. Academic Press, 1996.
- [31] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*, vol. 1. MIT press Cambridge, 1998.
- [32] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, no. 3-4, pp. 279–292, 1992.

# Improvements in speed and quality of learning in multi-agent systems using a novel expertness criteria and fuzzy integral

Dariush Hasanpour Adeg

d.hasanpour@ec.iut.ac.ir

Department of Electrical and Computer Engineering  
Isfahan University of Technology, Isfahan 84156-83111, Iran

Degree: M.Sc.

Language: Farsi

Supervisor: Assoc. Prof. Maziar Palhang (palhang@cc.iut.ac.ir)

## Abstract

In the real world, usually, peoples are coming together for sharing their knowledge and talking from their good and bad experiences and more or less everybody has something to say. Although we cannot ignore anybody's knowledge but it's common sense to assign more weight on the most experienced person's knowledge when we are going to decide what we need to do based on consultation from people. The achievements of this research have the same philosophy, that everybody needs to be heard. Fuzzy integrals are one of the most powerful and flexible methods for hearing everybody's knowledge and extract knowledge which is useful for everybody.

One of the challenges is that how to fairly answer the "what is the agents' expertise and how to determine the most and least expert agent?" question. To answer this question, in this thesis, we have proposed «the hypothesis of expertness» which defines a framework for "expertness criteria" definitions, and based on this framework we have introduced a new expertness criteria and showed that the defined framework and criteria are much more efficient than the state of the art criteria "Shortest Experienced Path" criteria. Also, the power of using fuzzy integrals for intelligence aggregation and non-additive measuring/knowledge is demonstrated.

## Key Words:

Multi-agent Systems, Cooperative Learning, Reinforcement Learning, Fuzzy Integral, Non-additive Knowledge



Isfahan University of Technology

Department of Electrical and Computer Engineering

# Improvements in speed and quality of learning in multi-agent systems using a novel expertness criteria and fuzzy integral

A Thesis

Submitted in partial fulfillment of the requirements  
for the degree of Master of Science

by

Dariush Hasanpour Adeh

Evaluated and Approved by the Thesis Committee, on Jan. 9. 2016

1. Maziar Palhang, Assoc. Prof. (Supervisor)
2. Abdolreza mirzaie, Asst. Prof. (Examiner)
3. Mohamad Hosein Manshaie, Asst. Prof. (Examiner)

Mohamad Reza Taban, Department Graduate Coordinator

