

## انتقال دانش میان عاملهای غیرهمسان به کمک اشتراک رویداد و اشتراک دانش

سید محمد رضا میرفتاح      مجید نیلی احمدآبادی      بابک نجار اعرابی

قطب علمی کنترل و پردازش هوشمند و آزمایشگاه هوش مصنوعی و رباتیک، گروه مهندسی برق و کامپیوتر، دانشکده فنی،  
دانشگاه تهران

پژوهشکده علوم شناختی، مرکز تحقیقات فیزیک نظری و ریاضیات

araabi@ut.ac.ir      mnili@ut.ac.ir      mirfataa@yahoo.com

**چکیده:** موضوع این نوشتار، ارائه روشهایی است که همکاری در یادگیری را میان عاملهای هوشمند غیرهمسان میسازند. پیش از این، مفید بودن همکاری در یادگیری میان عاملهای همسان به اثبات رسیده است. لیکن چگونگی همکاری در یادگیری میان عاملهای غیرهمسان که از لحاظ توانایی های عملکردی متفاوت هستند، به تازگی مورد توجه محققین علم هوش مصنوعی قرار گرفته است. در این مقاله، روشهایی برای همکاری مفید میان عاملهای غیرهمسان به هنگام یادگیری، ارائه شده است. به این ترتیب، عاملهایی که دارای اعمال مختلف هستند، می توانند با یکدیگر ارتباط معنی دار برقرار کرده و به نحو مطلوبی از دانش و تجارب دیگران برای ارتقاء دانش و کارایی خود استفاده کنند.

**کلمات کلیدی:** عاملهای غیرهمسان، آبرجدول نگاشت اعمال، همکاری در یادگیری تقویتی، اشتراک رویداد، اشتراک دانش

### ۱. مقدمه

اگر به چگونگی فرایند آموختن بیندیشیم، احتمالاً اولین ایده ای که به نظرمان می رسد، آموختن از طریق تعامل با محیط اطراف به همراه دریافت پاداش یا تنبیه از آن است. در این نوع یادگیری که یادگیری تقویتی<sup>۱</sup> نام دارد، عامل در جستجوی حالتی از محیط است که پاداشهای کسب شده را حداکثر و تنبیهها را حداقل کند. یکی از روشهای پیاده سازی این نوع یادگیری، Q-Learning می باشد. در این روش، دانش عامل توسط یک جدول نشان داده می شود. بازای هر حالت محیط، یک سطر و بازای هر عملی که عامل توانایی انجام آن را داشته باشد، یک ستون در این جدول قرار دارد. مقادیر موجود در خانه های این جدول نشان می دهند که در هر حالت، میزان مطلوبیت هر عمل چقدر است. هر چه این مقدار بیشتر باشد، عمل مربوطه مطلوبتر بوده و مسیر بهتری را برای رسیدن به هدف نشان می دهد. عامل با سعی و خطا و انجام دادن اعمال مختلف و گرفتن پاداش یا تنبیه از محیط، مقادیر موجود در جدول را بروز می کند.

تجربه کردن همه اعمال برای همه وضعیت ها، ممکن است از توان یک عامل خارج بوده یا لافل مستلزم صرف زمان طولانی برای او باشد. برای بالا بردن کارایی عاملها می توان از دانش کسب شده توسط دیگران نیز استفاده کرد. همکاری در یادگیری در جوامع انسانی و حیوانی دارای نمونه های فراوان است. در محیط عاملهای هوشمند نیز کارایی این روش نشان داده شده است (مثلاً در [۱] نشان داده شده که استفاده بجای از داده های بدست آمده از حواس دیگر عاملها و یا به اشتراک گذاشتن تجارب کسب شده توسط آنها موجب تسریع فرایند یادگیری می شود).

در بیشتر تحقیقات انجام شده روی سیستمهای چند عامله فرض بر همسان بودن عاملها قرار دارد. اما آنچه واقعاً در طبیعت مشاهده می شود، غیر از این است؛ حتی رباتهای تجاری به ظاهر یکسان نیز از لحاظ ویژگیهای حسی و حرکتی دقیقاً مانند هم نیستند. بعد از کار Ming Tan در سال ۱۹۹۳ [۱] که یادگیری تقویتی را بین عاملهای همسان بررسی کرد، در سال ۱۹۹۶ مقاله ای مشابه در مورد عاملهای غیر همسان منتشر شد [۲]. در این کار نویسندگان

<sup>۱</sup> Reinforcement learning

<sup>۲</sup> Homogeneous

با متفاوت انتخاب کردن نرخ یادگیری عاملها، آنها را غیرهمسان کرده‌اند. [۳] در سال ۲۰۰۰ به بررسی تقلید از یک عامل متفاوت در توانایی های فیزیکی (حرکتی) پرداخته است. در این کار عامل مقلد (آموزنده) با انجام تست امکان پذیری و تست جایگزینی، از مربی که عاملی است خبره، برای سریعتر پر کردن جدول Q خود کمک می‌گیرد. نتایج تحقیق دیگری که بدنبال پاسخ به این سؤال بوده که چگونه گروهی از عاملهای غیرهمسان، که در حل مسائل مشابهی دخیل هستند، می‌توانند به کمک مبادله مصلحت، همکاری کرده و کارایشان را افزایش دهند، در [۴] ارائه شده است. غیرهمسانی مورد نظر در این کار ناشی تفاوت در روشهای یادگیری عاملها بوده است. جدیدترین مقالات در زمینه همکاری در یادگیری میان عاملهای غیرهمسان از لحاظ توانایی های فیزیکی، به چگونگی استفاده از چارچوب ریاضی ارائه شده برای تطابق [۵]، در یادگیری اعمال به کمک تقلید پرداخته است. در این تحقیقات مکانیزم ALICE<sup>۱</sup> به عنوان روشی جامع برای مشخص کردن تطابقات با استفاده از هر روش ایجاد رفتار تقلیدی، معرفی شده است [۶]. خلاصه جامعی از تحقیقات مرتبط با موضوع این مقاله در [۷] ارائه شده است.

سؤال اصلی این تحقیق در مورد چگونگی انتقال و اشتراک دانش در گروهی از عاملهای یادگیر غیرهمسان<sup>۵</sup> است. غیرهمسانی این عاملها در اعمالی است که قابلیت اجرای آن را دارند. روش یادگیری عاملها، One-step Q-learning بوده و مسأله مورد نظر، توسط فرایند تصمیم گیری مارکوف<sup>۶</sup> مدل شده است. پیش از این [۱][۲]، انتقال دانش میان عاملهای همسان مورد بررسی قرار گرفته و روشهای خوبی نظیر اشتراک وزن دار استراتژی نیز برای این منظور طراحی شده است [۸]. لیکن صحت عملکرد تمام این روشها منوط به همسان بودن عاملها و یکسان بودن دانش بهینه آنها است. اگر عاملها از لحاظ مجموعه اعمالی که در اختیار دارند متفاوت باشند، دانش نهایی آنها متفاوت بوده و در نتیجه موضوع انتقال دانش با مشکل روبرو می‌شود. در واقع باید پاسخی برای چنین پرسشهایی یافت: آیا دانش یک عامل برای عامل غیرهمسان دیگر قابل استفاده است یا نه؟ آیا می‌توان با بکارگیری روشهایی، دانشهای غیرهمسان را به اشتراک گذاشت؟

در فصل بعد به تعریف مساله مورد نظر خود پرداخته و ضمن فصل بعد از آن، دو راه حل پیشنهادی خود را، ارائه می‌نمایم. در فصل ۴ به ذکر نتایج آزمایش‌هایی که نشان دهنده صحت عملکرد روش پیشنهادی است، پرداخته شده است. جمع بندی و بیان نتایج در فصل ۵ آمده است.

## ۲. تعریف مساله

مساله مورد بررسی در این مقاله قابلیت مدل شدن توسط فرایندهای تصمیم گیری مارکوف را دارد. در این مساله فرض شده که عاملها از هم مستقل بوده و عملکرد آنها با هم تداخل ندارد. در MDPهای مورد بررسی این تحقیق<sup>۷</sup>، فرضیات زیر در نظر گرفته شده‌اند:

الف) فضای حالت (S) برای تمام عامها یکسان است. ب) هر عامل از مجموعه اعمال خود (A) مطلع است. این اطلاع تنها به معنای دانستن تعداد اعمال و تمایز قائل شدن بین آنهاست. عامل از نتیجه عمل خود از قبل اطلاعی ندارد. پ) تابع انتقال حالت (T) برای عاملها مشخص نبوده و عاملها بدنبال یادگیری آن نیستند (یادگیری غیرمبتنی بر مدل). همچنین احتمال وجود اغتشاش در محیط وجود دارد. ت) تابع پاداش (R) برای عاملها معلوم نیست اما برای تمام عاملها یکسان است. این تابع، بصورت تابعی از حالت عامل (و نه حالت-عمل) تعریف شده است. مسأله مورد بررسی دارای یک هدف بوده و تنها با بازی رسیدن به آن، پاداشی (مثبت) به عاملها تعلق می‌گیرد. در سایر حالات، عامل، تنبیه (با مقادیر مختلف منفی) می‌شود (مثلاً بازی اتلاف انرژی، یا رفتن به حالات مخاطره آمیز). بستر آزمایشی که برای یادگیری در نظر گرفته شده، حل یک ماز<sup>۸</sup> است. در این مساله، هدف، پیدا کردن مسیر مناسبی از میان موانع، برای رساندن عامل به یک هدف از پیش تعیین شده می‌باشد. دو عامل غیرهمسان بطور مستقل برای یادگیری این ماز تلاش می‌کنند. یکی از این دو عامل در لحظات مختلف از دانش عامل دیگر استفاده می‌کند.

مساله غیرهمسانی در دو نوع مورد بررسی قرار گرفته است. در نوع اول، مجموعه اعمال هر دو عامل یکسان بوده و تنها ترتیب آنها متفاوت است. به عنوان مثال، مجموعه {U, D, R, L} را برای حرکت در چهار جهت اصلی در نظر گرفته، برای یک عامل ترتیب (U, D, L, R) و برای دیگری ترتیب (D, R, U, L) را در نظر می‌گیریم. در نوع دوم مجموعه اعمال یک عامل زیرمجموعه اعمال دیگری است. در ضمن تأثیر عمل (یا اعمال) اضافی، توسط جایگشتی از بعضی اعمال مشترک، قابل بازسازی است. به عنوان مثال، یکی دارای مجموعه اعمال پنج حرکتی {U, D, R, L, UR} بوده و دیگری دارای مجموعه اعمال {U, R, L, D} است. پیش از این در [۹][۱۰] چگونگی آموختن نگاهت بین اعمال در هر دو حالت غیرهمسانی

<sup>۲</sup> Correspondence

<sup>۳</sup> Action Learning for Imitation via Correspondence between Embodiments

<sup>۴</sup> Heterogeneous

<sup>۵</sup> ویژگی این مدل این است که هر انتقال حالت و سیگنال تقویتی که عامل با آن مواجه می‌شود، تنها به حالت فعلی و عملی که در آن انتخاب کند، بستگی داشته و به حالات قبلی ارتباطی ندارد.

<sup>۶</sup>  $M = \langle S, A, T, R \rangle$

<sup>۷</sup> Maze Problem

فوق، مورد بررسی قرار گرفته و برای این منظور استفاده از جدول نگاشت اعمال<sup>۹</sup> پیشنهاد شده است. اطلاعات کامل در مورد این روش در [۷] ارائه شده است.

نکته اصلی این مقاله، گسترش روشهای همکاری در یادگیری (من جمله WSS [۸]) به دامنه عملهای غیرهمسان است. در اینجا می‌خواهیم با استفاده از مقدماتی که فراهم شده، انتقال دانش بین عملهای غیرهمسان را بصورتی مفید و کارا مورد بررسی قرار دهیم. در این مقاله فرض شده که عملهای همکاری، نگاشت بین اعمال یکدیگر را در اختیار داشته و می‌خواهند برای بالا بردن سطح دانش خود، در حد امکان از دانش دیگران بهره‌برداری کنند. واضح است در صورتی که غیرهمسانی نوع اول به کمک داشتن AMT مناسب، رفع شده باشد، انواع روشهای همکاری در یادگیری که برای عملهای همسان پیشنهاد شده، برای چنین عملهای غیرهمسانی نیز قابل استفاده است. در واقع با بکارگیری AMT، می‌توان عملها را همسان در نظر گرفت. مسأله اصلی این مقاله در مورد تأثیر غیرهمسانی نوع دوم و در اصل حضور اعمال اضافی در مجموعه اعمال بعضی از عملهای همکاری است. در این مقاله فرض می‌کنیم عاملی که دارای عمل (یا اعمال) اضافی است، از عامل دیگر خبره‌تر است. در نتیجه، آنها را خبره و غیر خبره می‌نامیم. هدف این است که با طراحی روشهایی، با کمک دانش خبره، سرعت یادگیری غیرخبره را افزایش دهیم. در مورد فرض معکوس، یعنی خبره‌تر بودن عامل دیگر به [۷] مراجعه کنید.

### ۳. انتقال دانش

عملهایی که در این مقاله مورد توجه قرار دارند، از روش Q-Learning استفاده می‌کنند. در این روش غیرمبتنی بر مدل حل MDPها، بازی هر حالت و هر عمل، یک ارزش  $(Q(s, a))$  محاسبه می‌شود. این ارزش، مقدار مطلوبیت انتخاب عمل مربوطه را در آن حالت نشان می‌دهد. در تحقیق حاضر از نمایش جدولی مقادیر  $Q(s, a)$  استفاده شده است.

نکته‌ای که در اینجا اشاره به آن مورد نظر بوده، این است که در روش Q-Learning، داشتن اعمال اضافی باعث بزرگتر شدن (احتمالی) مقادیر ارزش  $Q(s, a)$ ها می‌شود. به تعبیر حساسی‌تر، اگر گزینه‌های جدیدی که در رسیدن به یک هدف، در اختیار انسان قرار می‌گیرند، هیچ مزیتی نسبت به گزینه‌های قبلی نداشته باشند، مسلماً در انجام عمل، از آنها استفاده نخواهد شد؛ یعنی اینکه استفاده از آنها وقتی صورت می‌گیرد که در رسیدن به هدف کمک مؤثری داشته باشند. برای اثبات مطلب فوق، به رابطه بروزسانی مقادیر Q توجه می‌کنیم:

$$Q(s, a) = (1 - \beta)Q(s, a) + \beta(r + \gamma V(t)) \quad \text{رابطه ۱-۳}$$

این رابطه بازی انجام عمل a در حالت s و رفتن به حالت t و دریافت سیگنال تقویتی  $\gamma$  روی  $Q(s, a)$  اعمال می‌شود. در این رابطه،  $\beta$  نرخ یادگیری و  $\gamma$  پارامتر کاهش اثر مطلوبیت حالات آتی بوده و  $V(t)$  از رابطه زیر بدست می‌آید:

$$V(t) = \max_b Q(t, b) \quad \text{رابطه ۲-۳}$$

عملگر max در رابطه اخیر، سبب می‌شود تا ارزش بهترین عمل در حالت t به حالت s تسری پیدا کند. در نتیجه اگر در حالتی، عمل اضافی مورد بحث، به عنوان بهترین عمل انتخاب شود، نسبت به عاملی که این عمل اضافی را ندارد، مقدار بزرگتری را برای انتشار در جدول Q ایجاد می‌کند. این مقدار بزرگتر می‌تواند در جاهای مختلفی از جدول Q تأثیر گذاشته و ارزش آنها را بالاتر ببرد. این، دقیقاً همان مشکلی است که عمل اضافی در روند همکاری در یادگیری طراحی شده برای عملهای همسان ایجاد می‌کند: دست بالا گرفتن ارزش بعضی عمل-حالت‌ها.

در روشهای معمول همکاری در یادگیری مانند WSS، فرض می‌شود که جدول بهینه دانش عاملها، یکسان است. به این ترتیب، ترکیب دانش عاملها باعث دست بالا گرفتن ارزش جایی نخواهد شد. اما با حضور عمل اضافی، احتمال غیریکسان شدن دانشهای بهینه و در نتیجه، مخدوش شدن دانش عاملها در اثر همکاری در یادگیری وجود دارد. بعنوان مثال، دانش بهینه دو عامل غیرهمسان از نوع دوم که در یک صفحه مشبک با مانع، برای رسیدن به هدفی واحد آموزش دیده‌اند، در جدول ۱-۳ دیده می‌شود. مقدار موجود در هر خانه، ارزش بهترین عمل را در آن خانه نشان می‌دهد.

فرض شده که  $\gamma = 0.5$  و پاداش رسیدن به هدف برابر ۱۰۲۴ و در غیر این صورت، صفر باشد. اگر دانش عامل دوم در اختیار عامل فاقد عمل اضافی قرار گیرد، عامل از خانه مشخص شده با دایره، با اجرای سیاست حریصانه، نمی‌تواند خارج شود؛ زیرا، ارزش این خانه بدلیل یک انتقال غیرممکن برای این عامل، بالا رفته است<sup>۱۰</sup>. می‌توان ایده‌هایی را که در این مقاله ارائه می‌شود، برگرفته از عبارات زیر دانست:

«اگر عمل اضافی در مسیر بهینه تا هدف، مورد استفاده قرار نگیرد باشد، دانش عامل در طی این مسیر، متأثر از وجود عمل اضافی نبوده و می‌توان آن را براحتی در اختیار عامل فاقد آن، قرار داد. در غیر این صورت، دانش عامل در طی مسیر تحت تأثیر عمل اضافی بوده و ممکن است موجب

<sup>۹</sup> Action Map Table (AMT)

<sup>۱۰</sup> نمونه‌ای از وابستگی دانش به ابزار کسب آن

سرردگمی عامل دیگر شود. بنابراین در چنین مواردی باید اثر این عمل را حذف یا جبران کرد.»

جدول ۱-۳: نمونه‌ای از محدودش شدن عملکرد عاملهای غیرهمسان در اثر همکاری در یادگیری

	۱۲۸	۱۰۲۴	هدف		۱۶	۱۰۲۴	هدف
$A_2 = A_1 \cup \{↗\}$	۲۵۶	۵۱۲	۱۰۲۴		۳۲	۶۴	۱۰۲۴
	۲۵۶	۲۵۶	۵۱۲	۵۱۲	۶۴	۱۲۸	۲۵۶

$$A_1 = \{\leftarrow, \rightarrow, \uparrow, \downarrow\}$$

### ۱-۳. اشتراک رویداد میان عاملهای غیرهمسان

هر رویداد تشکیل شده از رشته‌ای از سه‌تایی‌های (سیگنال تقویتی، عمل، حالت). این رشته از سه‌تایی‌ها با رسیدن به هدف و یا طی شدن تعداد مشخصی از گامهای جستجو در فضای مسأله خاتمه می‌یابد. عامل خبره (که دارای عمل اضافی است) در هر دور<sup>۱۱</sup> از عملکرد خود، رویدادی را تولید می‌کند. این رویداد به عامل غیرخبره داده می‌شود. عامل غیرخبره، این رویداد را از انتها به ابتدا (بصورت معکوس زمانی)، در ذهن خود تکرار می‌کند<sup>۱۲</sup>؛ گویی که خود او این رویداد را تجربه کرده است. فرض کنیم که رویداد مورد نظر بصورت زیر باشد:

$$\text{رابطه ۳-۳} \quad (s_1, a_1, r_1), \dots, (s_i, a_i, r_i), \dots, (s_n, a_n, r_n)$$

در نتیجه، برای بروز رسانی ارزش در حالت  $i$  خواهیم داشت:

$$\text{رابطه ۳-۴} \quad Q'(s_i, a'_i) = (1 - \beta)Q'(s_i, a'_i) + \beta(r_i + \gamma V'(s_{i+1}))$$

در نوشتن این روابط، از متغیرهای بدون اندیس برای عامل خبره استفاده شده است. اگر  $a_i$  یکی از اعمال مشترک باشد،  $a'_i$  عمل معادل در عامل غیرخبره برای آن است (که با توجه به AMT بدست می‌آید). اگر  $a_i$  عمل اضافی عامل خبره باشد، بررسی دو حالت لازم است: یا اینکه در حالت  $s_i$ ، هیچ یک از معادلهای  $a_i$  که در HAMT موجودند، قابل اعمال نیستند و یا اینکه لااقل یکی از آنها قابل اجراست. برای مثال در شکل ۱-۳، قسمتی از یک صفحه مشبک با مانع دیده می‌شود. در شکل سمت راست، عمل  $\rightarrow$  در حالت ۳ هیچ معادل قابل اجرایی ندارد. اما در شکل سمت چپ، می‌توان با اجرای  $\rightarrow$  و سپس  $\uparrow$  به نتیجه مشابه رسید.

در حالتی که هیچ یک از معادلهای ممکن نباشند، از سه‌تایی مربوطه صرف‌نظر کرده و به سه‌تایی بعدی پرداخته می‌شود. در حالتی که لااقل یکی از معادلهای قابل اجرا باشند، رابطه زیر برای بروز رسانی ارزش مربوطه بکار برده می‌شود:

$$\text{رابطه ۳-۵} \quad Q'(s_i, a'_i) = (1 - \beta)Q'(s_i, a'_i) + \beta(p + \gamma[r_i + \gamma V'(s_{i+1})])$$

در این رابطه  $a'_i$ ، اولین عمل از عمل ترکیبی معادل ممکن و  $p$  کوچکترین سیگنال تقویتی است.

		۱	
$A_2 = A_1 \cup \{↗\}$	۳	۲	

$$A_1 = \{\leftarrow, \rightarrow, \uparrow, \downarrow\}$$

### شکل ۱-۳: راست - معتبر نبودن معادلهای چپ - معتبر بودن یک معادل

صحت این رابطه را ضمن یک مثال بررسی می‌کنیم. در شکل ۲-۳ عامل خبره برای رفتن از حالت ۳ به حالت ۱، عمل  $\rightarrow$  را بکار برده است. از آنجایی که اجرای این عمل برای غیرخبره ممکن نیست، عامل غیرخبره می‌تواند با در نظر گرفتن معادل ترکیبی،  $\rightarrow$  را دنبال کند. با توجه به اینکه در مسأله مورد نظر تنها یک پاداش و آنهم به ازای رسیدن به هدف به عامل تعلق می‌گیرد، می‌توان مطمئن بود که رفتن به حالت ۲، با پاداش روبرو نخواهد بود<sup>۱۳</sup>. اما عامل غیرخبره از حالت ۲ بی‌اطلاع است (زیرا عمل  $\rightarrow$  را واقعاً اجرا نمی‌کند). در نتیجه نمی‌تواند از سیگنال تقویتی مربوط به آن هم مطلع باشد. راه حلی که در اینجا پیشنهاد می‌شود، بکار گرفتن کوچکترین سیگنال تقویتی ممکن برای حالت ۲ می‌باشد. ارزش حالت ۲ را نیز می‌توان با  $r_3 + \gamma V'(1)$  جایگزین کرد (توجه شود که در مسأله مورد نظر ما،  $I$  تابعی است از حالت جدید عامل).

<sup>۱۱</sup> Trial

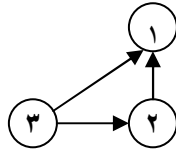
<sup>۱۲</sup> Mental replay

<sup>۱۳</sup> اگر رویداد تولید شده، بهینه نباشد، ممکن است که حالت ۲، همان هدف مسأله باشد. در این صورت باز هم تمییز شدن عامل باعث دست بالا گرفتن ارزش حالت ۳ نخواهد شد.

در نتیجه رابطه زیر برای بروز رسانی ارزش عمل  $\rightarrow$  در حالت ۳ مورد استفاده قرار می گیرد:

$$Q'(3, \rightarrow) = (1 - \beta)Q'(3, \rightarrow) + \beta(p + \gamma[r_3 + \gamma V'(1)]) \quad \text{رابطه ۳-۶}$$

تعمیم رابطه ۳-۵ به حالتی که معادل عمل اضافی، از ترکیب بیش از دو عمل حاصل شود، براحتی قابل انجام است.



شکل ۳-۲: انتقال ارزش به کمک معادل قابل اجرا

### ۲-۳. اشتراک دانش

روش اشتراک رویداد در جاهایی که رویدادهای غیر بهینه تولید شده باشند، قابل استفاده بوده و کارایی خود را نشان داده است. اما گاهی اوقات ترجیح می دهیم تا به جای آنکه مجدداً چرخ را اختراع کنیم، از آنچه وجود دارد، استفاده کنیم. دانش ایجاد شده در طی مدت جستجوی یک عامل در فضای مسئله، بسیار ارزشمندتر از تجارب مقطعی اوست. همانطور که قبلاً بیان شد، این دانش متأثر از ابزار کسب آن بوده و در نتیجه احتمالاً تحت تأثیر حضور اعمال اضافی، تغییراتی داشته است. اگر بدون انجام ملاحظاتی، روشی مثل WSS برای عاملهای غیرهمسان بکار برده شود، ممکن است که دانش عاملها از اعتبار خارج شود [۱۱].

اطلاع مورد استفاده در کنار جدول Q، در این تحقیق، مسیر حرکت عامل خبره در فضای مسئله است. این مسیر را می توان مستقلاً از روی مشاهده خبره و یا از روی رویداد مربوطه استخراج کرد. به عبارت دیگر جدول Q روی مسیر حرکت عامل خبره، برای انتقال دانش به عامل غیر خبره مورد نیاز است. فرض کنیم که رویداد مورد نظر بصورت رابطه ۳-۳ باشد. اگر این رویداد بصورت بهینه تولید شده باشد<sup>۱۴</sup>، بهترین جواب از این روش حاصل خواهد شد. در صورتی که رویداد، ناشی از حرکت غیربهینه در محیط باشد، از انتقال دانش معتبر (فاقد دست بالا گرفتن بعضی ارزشها) نمی توان مطمئن بود. این مطلب در آزمایشهای انجام شده نیز مشاهده شده است.

روند انتقال دانش از عامل خبره تولید کننده رویداد به عامل غیر خبره و فاقد عمل اضافی در جدول ۳-۲ ارائه شده است. این روند، حالت کلی انتقال دانش را مورد توجه قرار داده و حتی نحوه انتقال دانش در مورد رویدادهای غیربهینه نیز در آن لحاظ شده است. البته همانطور که بیان شد، نمی توان اطمینان داشت که دانشی که از رویداد غیربهینه منتقل می شود، بدون عیب باشد. روال کاهش ارزش در جدول ۳-۲ مشاهده می شود.

بیان این نکته ضروری است که، انتقال دانش همواره باید روی مسیر بهینه (با توجه به دانش عامل<sup>۱۵</sup>) صورت گیرد. انتقال دانش روی مسیر غیربهینه می تواند در بعضی شرایط موجب ایجاد دانش معیوب گردد. در واقع باید گفت که تضمینی برای ایجاد دانش غیرمعیوب در این حالت نیست؛ زیرا اعتبار روش فوق الذکر به درست شمردن اعمال اضافی مؤثر در دانش انتقالی و حذف تأثیرات آن است. حال اگر مسیر طی شده، غیربهینه باشد، دیگر نمی توان به مقدار شمارنده استناد کرد (ممکن است مقدار حقیقی شمارنده بیش از مقدار بدست آمده بوده و در نتیجه دانش دست بالا گرفته شده ای منتقل شود). نکته دیگری که در روال انتقال دانش از آن استفاده شده، مقایسه دانش جدید با دانش قبلی است. در شبیه سازیهای انجام شده، ارزش تمام حالات با کوچکترین مقدار ممکن آنها آغازین دهی شده است. در نتیجه، روند تغییرات ارزش هر عمل -حالت همواره صعودی خواهد بود. اگر دانش انتقالی مقداری کمتر از دانش موجود در ذهن عامل داشته باشد (به علت کمی دانش عامل پیشنهاد دهنده)، با مقایسه آنها از انتقال دانش کم ارزش جلوگیری شده و روند صعودی آن دچار خلل نمی شود.

### ۴. آزمایشها

آزمایشات بسیاری برای اطمینان از صحت و چگونگی عملکرد روشهای بیان شده در شرایط گوناگون، انجام شد. آنچه که در ادامه خواهد آمد، تنها به عنوان نمونه ای برای نشان دادن کارایی این روشها است.

<sup>۱۴</sup> این رویداد می بایست به هدف ختم شده باشد، در غیر این صورت تنها می توان از روش قبلی (اشتراک رویداد) استفاده کرد.

<sup>۱۵</sup> حتی اگر خبره نباشد.

جدول ۴-۱: روال کاهش ارزش  $v$  به اندازه counter مرتبه

```
double Discount(double v, int counter)
{
    for (int k=1; k < counter+1; k++)
        v =  $\gamma$ *v + p
    return v
}
```

جدول ۴-۲: روند انتقال دانش از عامل خبره دارای عمل اضافی به عامل فاقد آن



در این آزمایشات از دو عامل غیرهمسان با مجموعه اعمال موجود در جدول ۴-۳ استفاده شده است. این عاملها در صفحه مشبک دارای مانع<sup>۱۶</sup> قرار گرفته و با روش Q-learning به یادگیری پرداخته اند. هر دو عامل از روش انتخاب عمل  $\epsilon$ -greedy استفاده کرده اند. برای  $\gamma$  مقدار  $0.9$  در نظر گرفته شده است. در صورت دستیابی هر عامل به هدف (که با \* در تصاویر مربوطه مشخص شده)، محیط  $100$  امتیاز به عنوان پاداش به عامل داده و عامل با قرارگیری تصادفی در محلی دیگر، دور تازه ای را آغاز می کند<sup>۱۷</sup>. بازی برخورد به موانع، عامل با امتیاز  $10$  - تنبیه شده و در جای قبلی خود می ماند. رفتن به هر خانه دیگر، با  $1$  - امتیاز تنبیهی از طرف محیط پاسخ داده خواهد شد.

جدول ۴-۳: مجموعه اعمال عاملهای مورد آزمایش

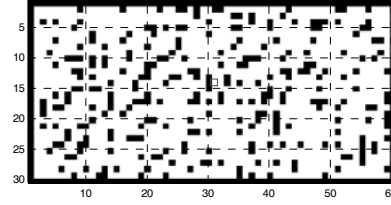
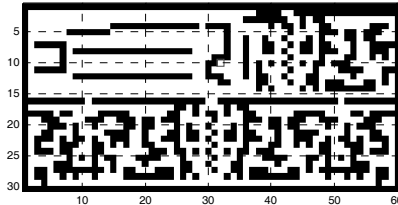
۵	۴	۳	۲	۱	شماره اندیس عمل
↗	↑	↓	→	←	مقصود عمل در عامل اول
	↓	↑	←	→	مقصود عمل در عامل دوم

۴-۱. کارایی روشهای اشتراک رویداد و انتقال دانش

برای نشان دادن کارایی روشهای مطرح شده، عاملها را در صفحه مشبک شکل ۴-۱ و شکل ۴-۲، قرار داده و نتیجه یادگیری فردی عامل دوم (بدون عمل اضافی) را با نتایج روشهای فوق که در آنها از عامل اول (که دارای عمل اضافی است) کمک گرفته شده، مقایسه می کنیم (هر نمودار میانگین  $10$

بار اجراست). شکل ۴-۱ با توجه به مجموعه اعمالی که عاملها در اختیار دارند، بعنوان صفحه ای که در آن معادلهای عمل اضافی، همواره قابل اجرا هستند، انتخاب شده است. شکل ۴-۲، بعنوان مسأله ای که در آن موانع محدود کننده (برای معادلهای عمل اضافی) وجود داشته و از پیچیدگی بالایی برخوردار است، انتخاب شده است. (آزمایشهای مشابهی بر روی صفحات مشبک دیگر انجام شده که شرح کامل آنها در [۷] آمده است).

در این آزمایشات، عامل اول از دانش بهینه برخوردار بوده است. برای به اشتراک گذاشتن رویداد، عامل اول با  $\epsilon=0.2$  در صفحه حرکت کرده و رویدادهای مشاهده شده را به عامل دوم منتقل می کند. نرخ یادگیری عامل دوم  $0.7$  می باشد. برای انتقال دانش، عامل اول با حرکت حریمانه ( $\epsilon=0$ ) روی صفحه، مسیرهایی را طی کرده و به اطلاع عامل دوم می رساند. برای یادگیری فردی، عامل دوم از نرخ یادگیری  $0.7$  و  $\epsilon=0.2$  استفاده کرده است (جدول ۴-۴). عامل دوم از HAMT صحیح مطلع است. از آنجا که در روشهای مورد بررسی، بروز رسانی به ترتیب عکس زمانی (از انتها به ابتدای مسیر) انجام می شود، نوع وارونه یادگیری فردی<sup>۱۸</sup> که در آن عامل دوم پس از اتمام هر دور یادگیری، بروز رسانی ارزشهای مسیر مربوطه را انجام می دهد، نیز مورد مقایسه قرار می گیرد.



شکل ۴-۱: صفحه مشبک ۶۰ در ۳۰، بدون مانع بد، هدف \* (Maze ۱) شکل ۴-۲: صفحه مشبک ۶۰ در ۳۰، دارای مانع بد و پیچیدگی زیاد، هدف \* (Maze ۲)

نقاطی که در آن ارزیابی های مربوط به همکاریها انجام شده، در دوره های ۱۰، ۵۰، ۱۰۰، ۲۰۰، ۳۰۰، ۴۰۰، ۵۰۰، ۶۰۰، ۱۰۰۰، ۲۰۰۰، ۷۰۰۰ و ۱۲۰۰۰ بوده است. ارزیابی دانش ناشی از یادگیریهای فردی در دوره های ۱۰، ۵۰، ۱۰۰، ۲۰۰، ۳۰۰، ۴۰۰، ۵۰۰، ۱۰۰۰، ۲۰۰۰، ۴۰۰۰، ۸۰۰۰ و ۱۲۰۰۰ انجام شده است. معیار اولی که برای مقایسه دانشها اطلاعات مفیدی را دربردارد، فاصله دانش عامل از مقدار بهینه آن بر حسب تعداد گامهای طی شده می باشد (رابطه ۴-۸).

از شکل ۴-۳ (که برای تمایز بهتر درمنحنیها، از مقیاسهای لگاریتمی برای هر دو محور آن استفاده شده) مطلب زیر استخراج می شود: روش اشتراک دانش به علت تأثیر گذاری روی مقدار ارزش بهترین عمل، و در نظر نگرفتن سایر اعمال، نمی تواند به دانش بهینه برسد. از آنجا که روش اشتراک رویداد محدود به حرکت بهینه در صفحه نیست، می تواند دانش مربوط به سایر اعمال (غیر از بهترین) را ارتقاء دهد.

$$\sum_s |Q'_{opt}(s, a) - Q'(s, a)| \quad \text{رابطه ۴-۷}$$

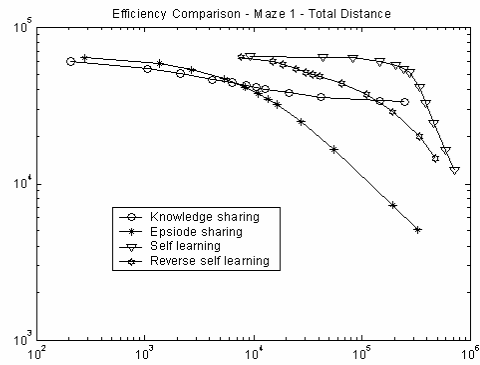
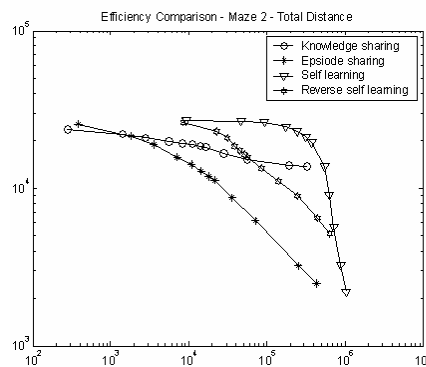
جدول ۴-۴: پارامترهای مربوط به سنجش کارایی روشهای انتقال دانش

Reverse Self Learning	Self Learning	Episode Sharing	Knowledge Sharing
$\alpha_2 = 0.7$	$\alpha_2 = 0.7$	$\alpha_2 = 0.7$	$\alpha_2 = 0.7$ (فقط در قسمت ES)
$\epsilon_2 = 0.2$	$\epsilon_2 = 0.2$	$\epsilon_1 = 0.2$	$\epsilon_1 = 0.0$

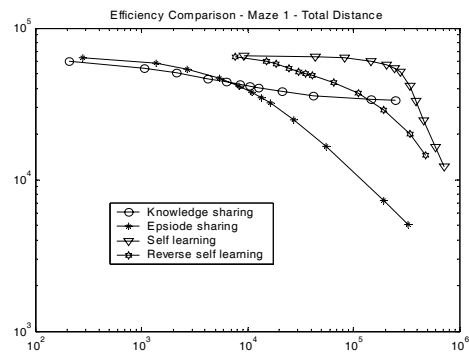
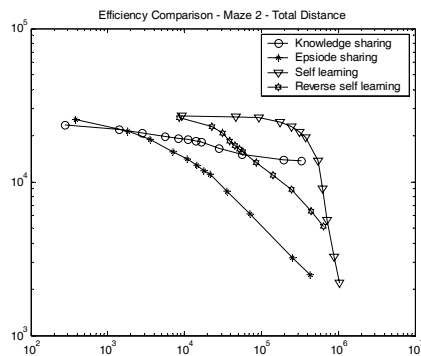
شکل ۴-۴ نشان دهنده نتیجه ارزیابی دانشها با استفاده از معیار رفتاری، بر حسب گامهای سپری شده در یادگیری است. برای محاسبه این معیار، حداقل فاصله هر خانه تا هدف بدست آمده و میانگین آنها مورد توجه قرار می گیرد. آنچه در شکل نشان داده شده، مقدار نرمال این میانگین با توجه به مقدار بهینه آن می باشد (محور افقی بصورت لگاریتمی نشان داده شده است). اگر نتوان از خانه ای به هدف رسید، فاصله آن، ۱۰۰۰ گام در نظر گرفته می شود.

از این شکل چند مطلب استخراج می شود: روشهای یادگیری فردی به مقدار بهینه خود خواهند رسید؛ هر چند که این کار نیاز به گذشت زمان طولانی تری دارد. بروز رسانی وارونه ارزش خانهها در مسیر، باعث تسریع در رسیدن به رفتار بهینه می شود. روشهای اشتراک دانش و اشتراک رویداد توانسته اند در ابتدای کار از یادگیری فردی پیشی بگیرند. به علت یک بودن تلوئیخی نرخ یادگیری در روش اشتراک دانش، این روش از اشتراک رویداد، سریعتر همگرا می شود. تفاوت مقدار نهایی روش اشتراک دانش و اشتراک رویداد، عمدتاً بدلیل تفاوت تعداد خانههایی است که از آنها می توان به هدف رسید (توجه کنید که هر خانه ای که از آن توان به هدف رسید، معادل ۱۰۰۰ گام خواهد بود). حداکثر دانشی که از روشهای اشتراک دانش و اشتراک رویداد حاصل می شود، لزوماً به اندازه بهینه خود نیست. علت این امر حضور موانع و تأثیر عمل اضافی در مسیرهایی است که به عنوان مسیر بهینه مشخص می شوند. در نتیجه، از بعضی از خانهها، حتی با حذف کردن اثر عمل اضافی، نمی توان بهترین مسیری را که در صورت نبود این عمل

بدست می آید را بدست آورد.



شکل ۴-۳: مقایسه فاصله ارزش بهترین عمل در هر حالت از مقدار بهینه آن بر حسب گامهای یادگیری



شکل ۴-۴: مقایسه فاصله دانش از مقدار بهینه آن بر حسب گامهای یادگیری

## ۲-۴. نتایج دیگر

یکی از نکاتی که در روش اشتراک دانش بسیار حائز اهمیت است، دریافت مسیر بهینه ناشی از دانش فعلی عامل کمک کننده است. همانطور که بیان شد با دور شدن این عامل از حرکت بهینه، احتمال ورود دانش معیوب به ذهن عامل دریافت کننده دانش بوجود می آید. صحت این مطلب در آزمایشهای انجام شده مشاهده شد. در این آزمایشها با مقایسه نتیجه حاصل از اشتراک دانش با بازی  $\epsilon$  های مختلف، مشاهده شد که دانشی که برای بهترین عمل منتقل می شود، در اثر افزایش  $\epsilon$ ، دچار عیب می شود. همانطور که انتظار می رود با افزایش میزان جستجو در صفحه، دانش کلی به نحو مطلوبتری کامل می شود، اما همین دانش نیز در اثر ادامه این روند، می تواند دچار عیب شود.

انتظار این است که با بیشتر شدن مقدار جستجو در محیط، روش اشتراک رویداد بتواند دانش بیشتری را منتقل کند. سرعت انتقال دانش در اثر افزایش  $\epsilon$ ، بیشتر می شود. البته واضح است که انتقال دانش روی بهترین عمل ضعیفتر می شود. در ضمن می توان دید که با افزایش  $\epsilon$ ، هر دور از گامهای بیشتری تشکیل می شود. مشاهده شد که معیار رفتاری دانش به تغییرات  $\epsilon$ ، حساسیت چندانی ندارد.

## ۵. جمع بندی و نتیجه گیری

نشان داده شد که بکارگیری روشهای موجود که برای انتقال دانش میان عاملهای همسان طراحی شده، در اینجا خیلی کارآمد نبوده و چه بسا مشکل آفرین باشد. بنابراین با توجه به دانستن نگاشت بین اعمال عاملهای غیرهمسان، سعی کردیم تا روشهایی را برای انتقال دانش بین عاملهای غیرهمسان (بویژه مسأله اصلی در نوع دوم غیرهمسانی خود نمایی می کند) بدست آوریم. نتیجه کار، دو روش اشتراک رویداد و اشتراک دانش بود.

دو روش ارائه شده در این مقاله، زمینه های کاربرد مخصوص به خود را دارند. در هر دو روش، انتقال دانش از عاملی خبره که دارای عمل اضافی است به عامل غیرخبره و فاقد آن مورد توجه بوده است. هنگامی که عامل فرستنده اطلاعات، بصورت حریصانه، فضای مسأله را طی می کند، استفاده از روش اشتراک دانش که وابسته به انتقال بهینه دانش عامل است، بسیار کاراست. در این روش با استفاده از دانش عامل خبره (دارای عمل اضافی) بعلاوه مسیری که در هر دور طی می کند، دانش عامل غیرخبره بروز می شود. در این روش، مسیری که عامل خبره بصورت حریصانه طی کرده از انتها به ابتدا مورد توجه قرار می گیرد. با شمارش تعداد اثرات عمل اضافی در هر نقطه از مسیر، دانش مربوطه برای استفاده عامل دیگر تطبیق داده می شود. از آنجا که



تعداد صحیح این اثرات در عمل تطبیق بسیار مؤثر است، لذا می‌بایست این مسیر، بهترین مسیری باشد که عامل دارای عمل اضافی با توجه به دانش خود آنرا انتخاب کرده باشد. در غیر این صورت، دانش منتقل شده ممکن است معیوب باشد. با استفاده از این روش، حداکثر دانش ممکن بین عاملهای غیرهمسان انتقال داده می‌شود. در صورتی که عامل خبره، عامل فاقد دارای عمل اضافی باشد، نمی‌تواند در مورد این عمل، دانشی را منتقل کند؛ مگر آنکه خود عامل دارای این عمل با انجام یک بازنگری در دانش خود، برای این عمل دانشی را ایجاد کند (این کار را می‌توان یکی از گامهای آتی در این تحقیق دانست).

اگر عامل خبره (دارای عمل اضافی)، فضای مسأله را بصورتی غیر بهینه جستجو کند، بعلاوه احتمال انتقال دانش معیوب، از روش فوق نمی‌توان استفاده کرد. در این شرایط، روش اشتراک رویداد که برای عاملهای غیرهمسان توسعه داده شده، نتایج خوبی را نشان داده است. در این روش، عامل غیرخبره با تکرار ذهنی رویداد تولید شده توسط عامل خبره بصورت از انتها به ابتدا، دانش خود را بروز می‌کند. این روش با توجه به نرخ یادگیری لحاظ شده در بروزرسانی دانش، از روش قبل کندتر عمل می‌کند. همچنین هر چه میزان جستجو در مسأله بیشتر باشد (دمای بالا در انتخاب عمل بولتزمن و  $\epsilon$  نزدیک به یک در روش  $\epsilon$ -greedy)، دانشی که منتقل می‌شود، گستردگی بیشتری روی فضای عمل - حالت دارد. این ویژگی در روش قبل که تنها منحصر به انتقال دانش مربوط به بهترین اعمال بود، وجود ندارد.

هر دو روش به نوع مسأله‌ای که در آن بکار گرفته می‌شوند، بسیار وابسته‌اند. اگر مسأله بگونه‌ای باشد که در جاهایی از آن عمل اضافی دارای معادل نباشد، انتقال دانش بخصوص از روش اشتراک دانش کم بهره خواهد بود. البته انتظار دیگری هم نمی‌توان داشت. اگر شخصی در مورد توانایی‌های منحصر به فرد خود، دانشی کسب کرده باشد، نمی‌توان انتظار داشت این دانش در شخصی که فاقد آن توانایی‌هاست، بکار آید. هر چه تأثیرات عمل اضافی در دانش تولید شده بیشتر باشد، دانشی که می‌تواند به اشتراک گذاشته شود، کمتر خواهد بود.

بهر حال توانستیم با استفاده از دانش عامل خبره و مسیر بهینه و یا با استفاده از رویداد تجربه شده توسط عامل خبره، دانشی را به عامل غیرخبره منتقل کنیم. این انتقال دانش از جهاتی بر روشهای یادگیری فردی (عادی و وارونه) برتری خود را نشان داد. به نظر می‌رسد که انتقال و تطبیق دانشهای غیرهمسان بدون استفاده از اطلاعی دیگر (مثل مسیر حرکت) کار ممکن نباشد. اما این مسأله بعنوان افقی از این کار می‌تواند موضوع تحقیقات دیگری باشد (هر چند که جستجو برای یافتن روشی که تنها از روی خود دانش بتواند، دانش منطبق با عامل غیرهمسان را مشخص کند، در این تحقیق به ثمر نرسید).

## ۶. مراجع

- [۱] Ming Tan, Multi-agent reinforcement learning: Independent vs. cooperative agents, in *Proceedings of the Tenth International Conference on Machine Learning*, ۱۹۹۳.
- [۲] Kawaishi, I., et al, Experimental Comparison of a Heterogeneous Learning Multi-Agent System with a Homogeneous one, *IEEE International Conference on Systems, Man and Cybernetics*, ۱۹۹۶, pp. ۶۱۳-۶۱۸.
- [۳] Bob Price & Craig Boutilier, Imitation and Reinforcement Learning in Agents with Heterogeneous Actions, *Proceedings of the Seventeenth International Conference on Machine Learning*, ۲۰۰۰.
- [۴] Luís Nunes & Eugénio Oliveira, Advice Exchange in Heterogeneous Groups of Learning Agents, *Tech. report 1-02 FEUP/ISCTE/LIACC-NIAD&R*, ۲۰۰۲.
- [۵] Nehaniv, C. L. & Dautenhahn, K., The Correspondence Problem, in Nehaniv, C. L. & Dautenhahn, K., editors, *'Imitation in Animals and Artifacts'*, MIT Press, ۲۰۰۲.
- [۶] Alissandrakis, A., Nehaniv, C. L. & Dautenhahn, K., Solving the Correspondence Problem Between Dissimilarly Embodied Robotic Arms Using the ALICE Imitation Mechanism, *Proceedings of the Second International Symposium on Imitation in Animals & Artifacts*, The Society for the Study of Artificial Intelligence and Simulation of Behavior, ۲۰۰۳, pp. ۷۹-۹۲.
- [۷] سید محمد رضا میرفتاح، همکاری در یادگیری میان عاملهای غیرهمسان، پایان نامه کارشناسی ارشد، دانشگاه تهران، ۱۳۸۲.
- [۸] Majid Nili Ahmadabadi & Masoud Asadpour, Expertness Based Cooperative Q-learning, *IEEE Transactions on Systems, Man and Cybernetics*, Part B: Cybernetics, Vol. ۳۲, No. ۱, ۲۰۰۲, pp. ۶۶-۷۶.
- [۹] S.M.Reza Mirfattah and Majid Nili Ahmadabadi, Cooperative Q-learning with heterogeneity in actions, *Proceedings of 2002 IEEE International Conference on Systems, Man and Cybernetics*, Hammamet, Tunisia, ۲۰۰۲.
- [۱۰] سید محمد رضا میرفتاح، مجید نیلی احمدآبادی و بابک نجار اعرابی، معادل یابی اعمال در عاملهای غیر همسان: گامی به سوی همکاری در یادگیری، مجموعه مقالات هشتمین کنفرانس بین المللی سالانه انجمن کامپیوتر ایران، ۱۳۸۱، صفحه ۱۳۴-۱۴۱.
- [۱۱] Mastour Eshgh, S., Araabi, B., Nili Ahmadabadi, M., Cooperative Q-learning through State Transitions: A Method for Cooperation Based on Area of Expertise, *Proceedings of the 4th Asia-Pacific Conference on Simulated Evolution and Learning (SEAL'02)*, Singapore, ۲۰۰۲.