



## مجرد سازی دانش و انتقال آن با استفاده از یادگیری تقویتی

مریم سادات میریان<sup>۲</sup>  
mmirian@ut.ac.ir

احمدرضا ولی<sup>۱</sup>  
ar.vali@gmail.com

مجید نیلی احمدآبادی<sup>۲</sup>  
mnili@ut.ac.ir

نرجس زارع<sup>۱</sup>  
zare.narjes@gmail.com

۱- دانشگاه صنعتی مالک اشتر، مجتمع برق و الکترونیک ۲- دانشگاه تهران، دانشکده برق و کامپیوتر، آزمایشگاه رباتیک

چکیده - هنوز فاصله قابل توجهی بین قابلیت یادگیری در سیستم های هوشمند و سیستم های بیولوژیکی دیده می شود. یکی از دلایل این مساله استفاده نکردن از دانش بدست آمده در طول یادگیری در یک سیستم هوشمند می باشد. به منظور نشان دادن کار آمدی این دیدگاه در این مقاله روشی برای انتقال یادگیری ارائه شده است. در این روش دانش یاد گرفته شده در کار مرجع بصورت مفاهیم مجرد در آمده و به کار هدف منتقل می شود. عامل با استفاده از این دانش منتقل شده می تواند به یادگیری سریعتری در کار هدف دست یابد. نتایج شبیه سازی نشان می دهند که روش ارائه شده باعث افزایش پاداش بدست آمده در طول یادگیری، بخصوص اوایل یادگیری و در نتیجه افزایش سرعت یادگیری می شود.

کلید واژه- انتقال یادگیری، مجرد کردن دانش، مفهوم، یادگیری تقویتی

### ۱- مقدمه

موفق خواهد بود که بعد از استفاده از دانش یاد گرفته شده از کار مرجع، یادگیری در کار هدف سریعتر و با عملکرد بهتری صورت گیرد. اکثر روش هایی که تا کنون در راستای انتقال دانش ارائه شده اند از نوعی نگاشت بین کارهای مرجع و هدف استفاده کرده اند [2]. یافتن چینی نگاشتی بین کارها ملزم داشتن اطلاعات زیادی در مورد کارها و همچنین وجود و شناسایی تشابهات بین کارها می باشد که این امر باعث ایجاد محدودیت در انتقال دانش می شود و کار را مشکل می سازد. برای حل این مشکل می توان از نگاشت در سطح بالایی از تجرید بین دو مساله استفاده کرد. در حالت کلی می توان گفت که مقوله های انتقال یادگیری و انتزاعی کردن حالت خیلی به هم وابسته و نزدیک بهم هستند [3]. در انتقال یادگیری عامل می کوشد که از دانش کار مرجع در کار هدف استفاده نماید. رسیدن به این هدف نیازمند این است که یک صورت انتزاعی از فضای حالت داشته باشیم تا بتوان دانش کار مرجع را در کارهای دیگر که حوزه های کاری متفاوت دارند، استفاده کرد. بنابراین مساله تصمیم گیری در مورد انتخاب دانش برای انتقال در حوزه های متفاوت

موجودات هوشمند با داشتن هوش طبیعی در تعامل با محیط یاد می گیرند که در برابر هر شرایطی از محیط چگونه رفتار کنند. در شاخه های مختلف هوش مصنوعی نیز سعی شده است که با الهام گرفتن از نحوه تعامل سیستم های هوشمند طبیعی با محیط روش هایی برای توسعه سیستم های هوشمند مصنوعی ارائه شود. یادگیری تقویتی یکی از اما هنوز فاصله قابل توجهی بین قابلیت یادگیری در سیستم های هوشمند و سیستم های بیولوژیکی دیده می شود. یکی از دلایل این مساله استفاده نکردن از دانش بدست آمده در طول یادگیری در یک سیستم هوشمند می باشد. به این منظور، اخیراً محققان توجه خود را به انتقال دانش در سیستم های هوشمند معطوف کرده اند.

هدف از انتقال دانش، یادگیری سریع کار (هدف) بعد از یادگیری کار (مرجع) متفاوت اما مشابه (وابسته) به آن می باشد [1]. انتقال دانش به عامل اجازه می دهد که ابتدا یک کار مرجع ساده اولیه را یاد بگیرد و سپس با توجه به آن کار پیچیده تری را یاد بگیرد. در صورتی انتقال یادگیری



فراموشی است. یک راه حل مرسوم، تخمین زدن تابع حالت و عمل بهینه یا همان تابع  $Q$  که حالات و اعمال را به ماکزیمم پاداش مورد انتظار که با شروع از حالت  $S$  و عمل  $a$  نگاهت می دهد، می باشد.

### ۳- انتقال یادگیری

انتقال یادگیری به معنی بکار بردن دانش یادگرفته شده در یک کار ( کار مرجع) به منظور بهبود یادگیری در کار دیگر (کار هدف) می باشد. بشر بطور قابل ملاحظه ای از دانش-های یادگرفته شده در کارهای گذشته اش برای یادگیری بهتر و سریعتر در کارهای خود بهره می برد. بیشتر روش هایی که تا کنون در راستای انتقال دانش ارائه شده اند از نوعی نگاهت بین کارهای مرجع و هدف استفاده کرده اند. یافتن چینی نگاشتی بین کارها ملزم داشتن اطلاعات زیادی در مورد کارها و همچنین وجود و شناسایی تشابهات بین کارها می باشد که این امر باعث ایجاد محدودیت در انتقال دانش می شود و کار را مشکل می سازد. برای حل این مشکل می توان از نگاهت در سطح بالایی از تجرید بین دو مساله استفاده کرد. در این مقاله برای نیل به این هدف از مجرد کردن دانش از طریق یادگیری مفاهیم استفاده شده است.

### ۴- مجرد کردن دانش

محیط واقعی که ما انسانها در آن زندگی می کنیم سرشار از اطلاعات پیوسته و گسسته که با نا یقینی و نویز همراه است می باشد. یادگیری در چنین محیط پیچیده ای، بدون استفاده از مکانیزم های خاصی غیر ممکن می نماید. یکی از مهم ترین حربه هایی که انسان برای برخورد با چنین محیط هایی بکار می برد مجرد کردن می باشد. مجرد سازی مکانیزمی است که طی آن فضای ادراکی پیچیده عامل هوشمند به یک فضای ساده تر که توسط عامل قابل مدیریت کردن می باشد نگاهت می شود و از آنجا که سعی می شود در این فرآیند تا حد امکان محتوای اطلاعاتی فضای ادراکی اصلی حفظ شود، فضای ادراکی حاصل به نوعی، باز نمایی سطح بالاتر و یا مجردتر از فضای اصلی خواهد بود [۶].

می تواند به مساله انتزاعی کردن حالت برای یک مجموعه از حوزه های مرجع تبدیل شود. در این مقاله سعی بر این است تا با استفاده از مفاهیم سلسله مراتبی، عامل دانش خود را برای استفاده مجدد در کار دیگر ( معمولاً پیچیده تر ) مجرد کند. رویکرد ما برای تعریف مفهوم و مدل کردن آن در ذهن عامل یک رویکرد کارکردی است [۴]. یعنی مفاهیم در فضای  $Q$ -value ها شکل می گیرند. نتایج بدست آمده نشان می دهد که روش پیشنهادی برای انتقال دانش نه تنها باعث سرعت در یادگیری می شود بلکه حتی با تغییر سنسورهای ربات ( فضای حسی) عامل می تواند از دانش قبلی به علت قابلیت تعمیم آن، استفاده کند و بوسیله آن سرعت یادگیری خود را افزایش دهد.

ترتیب مطالب این مقاله بدین صورت می باشد: در قسمت دوم مروری خلاصه بر یادگیری تقویتی داریم. قسمت سوم انتقال دانش و راهکارهای موجود و هدف از آن را بررسی می کند. بعد از آن در قسمت چهارم تعریفی از مجرد کردن دانش خواهیم داشت. سپس در قسمت پنجم به نحوه ایجاد مفاهیم و یادگیری آنها می پردازیم. در قسمت ششم به چگونگی انتقال یادگیری در این مقاله می پردازیم و در آخرین بخش نتایج حاصل از شبیه سازی آورده شده است.

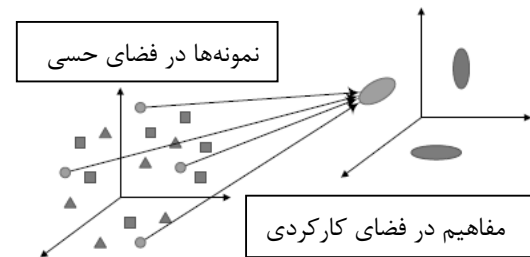
### ۲- یادگیری تقویتی

در چارچوب یادگیری تقویتی، عامل یادگیرنده در طی عمر خود  $t = 0, 1, 2, 3, \dots$  در تقابل با محیط می باشد. در هر مرحله زمانی  $t$  عامل حالت جاری محیط  $S_t$  را مشاهده می کند و بر اساس آن عمل  $a$  را انجام می دهد که باعث می شود محیط به حالت بعدی  $S_{t+1}$  منتقل شود و در نتیجه آن عامل پاداش  $r$  از محیط دریافت کند [۵]. در یک سیستم مارکف حالت بعدی محیط و پاداش دریافتی تنها به عمل و حالت قبلی عامل در محیط بستگی دارد. هدف عامل در یادگیری ماکزیمم کردن پاداش بدست آمده در طول زمان می باشد. عامل با یادگرفتن نگاهت حالات به اعمال که سیاست نامیده می شود، این کار را انجام می دهد. عبارتی هدف عامل انتخاب عمل بطوریکه مقدار پاداش مورد انتظار  $E[\sum_{t=0}^{\infty} \gamma^t r_{t+1}]$  را افزایش دهد، می باشد، که  $\gamma$  ضریب



## ۵- یادگیری مفاهیم

یکی از روشهای موجود برای مجردسازی، تقسیم فضای ادراکی عامل به یک سری کلاس های شباهت است بطوری- که هر کلاس حالت های مشابه در فضای ادراکی را در خود جای دهد. به هر یک از این کلاس های شباهت مفهوم گفته می شود. انسانها در تعامل با محیط یادگرفته اند که برای درک محیط پیرامونشان تنها به ویژگی های حسی اتکا نداشته باشند و با عمل هایشان محیط را بشناسند. آنها چون قابلیت استخراج مفاهیم مجرد از محیط را دارند، می-توانند یادگیری در یک محیط را به محیط های جدید تعمیم دهند. مفاهیمی که با اتکا به ارزش اعمال بدست می آیند را مفاهیم کارکردی می گویند. در این مقاله نشان داده شده است که مفاهیم کارکردی ابزار مناسبی برای مجرد کردن دانش ربات بشمار می رود که می تواند قابلیت مناسبی برای به اشتراک گذاشتن دانش بین عامل هایی که فضای حسی متفاوتی دارند یا بین دو مساله که فضای حسی متفاوتی دارند محسوب شود. شکل (۱) شماتیکی از نمونه ها در فضای حسی و مفاهیم استخراجی از آنها در فضای کارکردی نشان می دهد.



شکل ۱: شماتیکی از نمونه ها در فضای حسی و مفاهیم استخراجی از آنها در فضای کارکردی [۴]

## ۶- مجرد کردن دانش و انتقال آن با استفاده از

### یادگیری تقویتی

مجرد کردن و قابلیت تعمیم دو ویژگی مهم برای سیستم های هوشمند می باشد که سرعت و کیفیت یادگیری را بهبود می بخشد. در این مقاله برای مجرد کردن دانش،

دانش بدست آمده از یادگیری عامل در کار مرجع را بصورت مفاهیم مجرد در می آوریم. برای بدست آوردن مفاهیم، بردارهای  $Q$  بدست آمده عامل پس از یادگیری در مساله مرجع توسط روش خوشه بندی  $k$ -means بصورت دسته هایی مجزا در می آیند مراکز این دسته ها همان مفاهیم انتقالی هستند که عامل در مساله هدف از آنها برای یادگیری سریعتر استفاده می کند. پس از محاسبه مفاهیم آنها را به مساله هدف منتقل می کنیم. الگوریتم یادگیری در هر دو مساله مرجع و هدف یادگیری  $Q$  می باشد. عامل در کار هدف ابتدا یادگیری عادی خود را شروع می کند. سپس پس از اینکه به دانش لازم برای تشخیص صحیح مفهوم دست یافت از مفاهیم انتقالی برای تصمیم گیری استفاده می کند. معیاری که برای این کار در نظر گرفته شده است تعداد تکرار قرار گرفتن عامل در آن حالت می باشد. اگر عامل  $N_c$  بار یک حالت را مشاهده کرده باشد مجاز به استفاده از مفاهیم برای انتخاب عمل مناسب می باشد که این مقدار با سعی و خطا بدست می آید. این کار را چون عامل در ابتدای یادگیری دید کافی به مساله ندارد و نمی-تواند تشخیص درستی از مفهومی که در آن قرار دارد داشته باشد، انجام می دهیم. عامل پس از اینکه توانایی استفاده از مفاهیم را بدست آورد با توجه به بردار ارزش  $Q$  که در آن قرار دارد شبیه ترین مفهوم را که نزدیکترین مفهوم به بردار  $Q$  می باشد را انتخاب می کند. فاصله اقلیدسی در اینجا بعنوان شاخص شباهت در نظر گرفته شده است:

$$d_i = \sqrt{(\bar{Q}_s - \bar{Q}_{c_i})^T (\bar{Q}_s - \bar{Q}_{c_i})} \quad (1)$$

$$Q_c = \arg \min_{Q_c} (d) \quad (2)$$

اگر فاصله بردار  $Q$  تا نزدیکترین مفهوم از محدوده تعیین شده کم تر باشد. عامل با توجه به آن مفهوم تصمیم گیری می کند در غیر اینصورت با استفاده از بردار  $Q$  خود و الگوریتم یادگیری  $Q$  عمل مناسب را انتخاب می کند. با این کار عامل تنها در حالاتی که دانش مرجع برایش مفید است از دانش انتقالی استفاده می کند. با انجام عمل و انتقال محیط به حالت بعدی و پاداشی که محیط به عامل می دهد



شکار را متوقف کند. در صورتیکه شکارچی به مانع برخورد کند پاداش ۴- به او داده می شود. هزینه ای که شکارچی برای تیرهای بدون نتیجه می پردازد ۶- می باشد. این هزینه باعث می شود که شکارچی فقط در صورت نیاز ( در تیر رس بودن شکار ) تیراندازی کند. ۲۰+ امتیاز پاداش را شکارچی وقتی که شکار را بدون تیر زدن متوقف می کند ( در یکی از خانه های مجاور شکار قرار می گیرد ) دریافت می کند. اگر شکارچی با تیر شکار را متوقف کند پاداش ۳۰+ می گیرد. هر دوره یادگیری با متوقف شدن شکار پایان می یابد. پس از آن دوباره شکار و شکارچی بطور تصادفی در محیط قرار داده می شوند و دوره ای جدیدی آغاز می شود. این کار تا تمام شدن تعداد تکرارهای در نظر گرفته شده برای یادگیری ادامه می یابد. پارامترهای یادگیری بصورت زیر تنظیم می شوند:

نرخ یادگیری ( $\eta$ ) بطور نزولی با افزایش تعداد تکرارها از مقدار ۰.۸ تا ۰.۲ کاهش می یابد.

ضریب فراموشی ( $\gamma$ ) برابر با ۰.۹ در نظر گرفته شد. همانطور که قبلا گفته شد دانش بدست آمده عامل پس از یادگیری در کار مرجع ابتدا بصورت مجرد در آمده و سپس به کار هدف منتقل می شود. برای مجرد کردن دانش و در آوردن مفاهیم بردارهای ارزش  $Q$  نرمالیزه شده را با استفاده از روش خوشه بندی  $k$ -means دسته بندی می کنیم. مراکز خوشه ها که مفاهیم انتقالی هستند در این مساله به صورت زیر می باشند:

concepts	Value of 1th action	Value of 2th action(go right)	Value of 3th action(go left)	Value of 4th action(go down)	Value of fifth action(shoot)
$C_1$	0.064	0.069	0.063	<u>0.989</u>	0.0156
$C_2$	0.079	0.076	0.083	0.076	<u>0.9945</u>
$C_3$	0.067	0.069	<u>0.978</u>	0.062	0.0145
$C_4$	<u>0.987</u>	0.068	0.077	0.061	0.0173
$C_5$	0.071	<u>0.965</u>	0.066	0.064	0.0168
$C_6$	0.750	0.748	0.755	0.742	0.0802

همانطور که ملاحظه می شود ۶ مفهوم از محیط بدست آمده است. با توجه به مراکز بدست آمده، در هر بعد تنها یکی از

عامل با استفاده از یادگیری  $Q$  مقدار بردار  $Q$  خود را با استفاده از رابطه (۳) بروز می کند.

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a')) \quad (3)$$

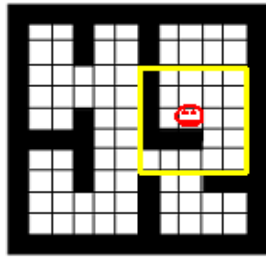
## ۷- نتایج شبیه سازی

مساله ای که برای شبیه سازی در نظر گرفته شده، مساله شکار و شکارچی می باشد. مساله شکار و شکارچی یکی از مسائل کلاسیک برای مطالعه و مقایسه روش های متفاوت یادگیری در هوش مصنوعی می باشد [۷].



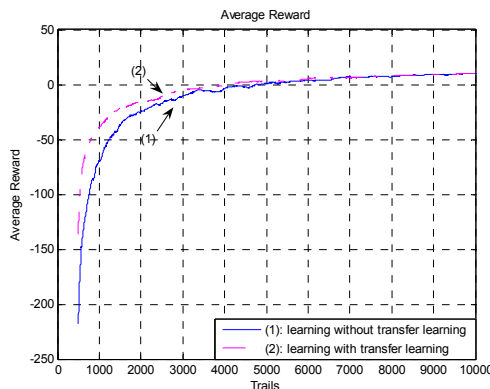
شکل ۲: نمایی از دو محیط شبیه سازی، شکل سمت چپ محیط شبیه سازی برای مساله مرجع و سمت راست برای مساله هدف شکار و شکارچی در یک محیط دو بعدی مربعی شکل که هر بعد آن به ۱۲ قسمت تقسیم شده است، واقع شده اند. دورتا دور محیط دیوار قرار گرفته است و تعدادی مانع در آن قرار داده شده است. عامل شکارچی در هر قدم یکی از پنج عمل بالا، پائین، راست، چپ و یا تیر زدن را با احتمالی که با تابع بولتزمن مشخص می شود انتخاب می کند. شکارچی بایستی بدون برخورد به مانع و در کمترین زمان ممکن شکار را بگیرد یا آن را با تیر بزند. برای گرفتن شکار کافی است که شکارچی در یکی از خانه های مجاور شکار قرار بگیرد. الگوریتم یادگیری شکارچی یادگیری  $Q$  می باشد و شکار هم با استفاده از روش میدان پتانسیل از شکارچی فرار می کند. در هر گام یادگیری پس از اینکه عامل یادگیرنده ( شکارچی ) عمل حاصل از تصمیم گیری خود را انجام داد بر اساس اینکه به چه ناحیه ای از محیط منتقل شده است از محیط پاداش دریافت می کند که این تابع پاداش بدین صورت می باشد:

بازای هر حرکت بدون نتیجه پاداش ۱- به شکارچی داده می شود. اینکار باعث می شود که شکارچی با کمترین حرکت



شکل ۴: ربات با سنسورهای جدید در محیط یادگیری

اما با وجود اینکه سنسورهای دو کار تغییر کرده باز هم مفاهیم بدست آمده از کار مرجع یادگیری ربات را بهبود بخشیده است. شکل (۵) نمودار پاداش متوسط را نشان می-دهد.

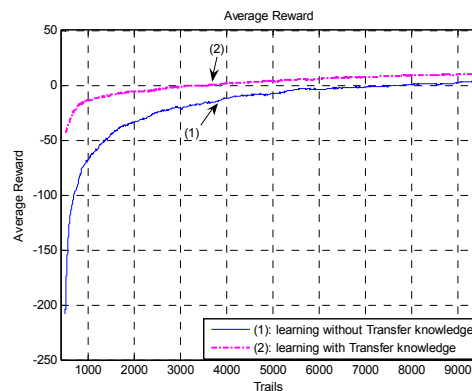


شکل ۳- مقایسه نمودارهای پاداش متوسط: (۱) یادگیری بدون انتقال دانش و (۲) یادگیری با انتقال دانش در دو مساله با فضای حسی متفاوت

در آخرین قسمت از شبیه سازی های انجام شده انتقال یادگیری بین یک محیط مارکف (MDP) و یک محیط مشاهده پذیر جزئی مارکف (POMDP) بررسی می شود. در یک محیط POMDP بعضی از حالت های متفاوت محیط از نظر عامل یکی می شوند و این امر باعث می شود که عامل به یادگیری مناسبی دست پیدا نکند. در این جا ما با انتقال دانش کسب شده از کار مرجع که MDP بوده باعث بهبود یادگیری عامل در محیطی که برای عامل POMDP است، می شویم. شکل (۶) حاکی از این امر می باشد.

اعمال که زیر آنها خط کشیده شده است از احتمال بالایی برای انتخاب برخوردارند.

با توجه به بردارهای مفاهیم بدست آمده می بینیم که هر مفهوم یک عمل خاص را توصیه می کند مثلا مفهوم اول عمل چهارم که در اینجا پایین آمدن است را توصیه می کند مفهوم دوم که نشانگر در تیررس بودن شکار می باشد عمل تیر زدن را توصیه می کند. عامل در مساله هدف پس از طی یک دوره یادگیری، در حالت هایی که بیش از ۵ بار در آنها قرار گرفته از دانش منتقل شده برای تصمیم گیری کمک می گیرد. نتایج حاصل از شبیه سازی در شکل (۳) آمده است. لازم به توضیح است که نمودارها، حاصل متوسط-گیری از ۵ بار اجرا هستند تا اطلاعات بدست آمده قابل اعتمادتر شوند و اثر حالت های خاص کم تر شوند.



شکل ۳- مقایسه نمودارهای پاداش متوسط: (۱) یادگیری بدون انتقال دانش و (۲) یادگیری با انتقال دانش

برای نشان دادن قابلیت تعمیم مفاهیم انتقالی و اینکه دانش انتقالی به فضای حسی عامل بستگی ندارد، سنسورهای ربات را در کار هدف عوض کرده و از همان دانشی که ربات با استفاده از سنسورهای قبلی خود بدست آورده به یادگیری در کار هدف می پردازیم. سنسورهای اولیه ربات موقعیت X و Y ربات شکارچی و X و Y شکار را می دهند. در صورتیکه سنسورهای ثانویه ربات همانطور که در شکل (۴) مشاهده می شود، تنها موقعیت خانه های اطراف ربات تا شعاع دو خانه مجاور را به ربات می دهد.

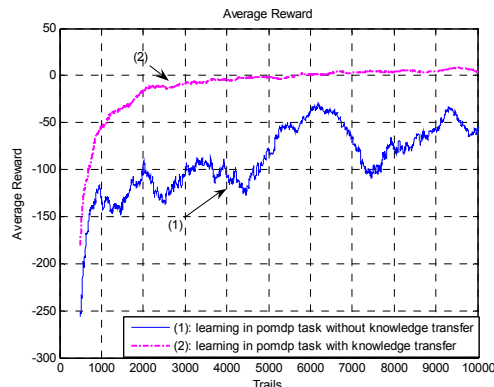
[3] Thomas J. Wash and Michael L. Litman. Transferring state Abstraction between MDPs. In ICML-07 conference.

[4] Hamide Vosoughpour, Majid Nili Ahmadabadi, Maryam S. Mirian, Babak Nadjar Araabi. Hierarchical Functional Concept Formation using Reinforcement Learning. In ICTA2007 conference.

[5] R. S. Sutton and A. G. Barto. Reinforcement Learning: An Introduction. MIT Press, Cambridge, MA, 1998.

[6]. Hadi Firouzi, Majid Nili Ahmadabadi, Babak N. Araabi. A Probabilistic Reinforcement-Based Approach to Conceptualization. In International Journal of Intelligent Technology (IJIT), Volume 3, pp 48-55, 2008

[7] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in Proc. Tenth Int. Conf. Machine Learning, Amherst, MA, June 1993.



شکل ۶- مقایسه نمودارهای پاداش متوسط: (۱) یادگیری بدون انتقال دانش و (۲) یادگیری با انتقال دانش بین دو مساله MDP و POMDP

## ۸- نتیجه گیری

در این مقاله به بررسی انتقال دانش بین دو کار در حالت-های مختلف پرداخته شد و نشان داده شد که اگر دانش بدست آمده از کار اول را بصورت مجرد در آوریم براحتی می توانیم آن را در یک کار مشابه دیگر استفاده کنیم حتی اگر فضای حسی دو کار با یکدیگر متفاوت باشد. در این روش دیگر نیازی به نگاشت یک به یک بین حالت و عمل عامل در دو کار نیست و مشکل پیدا کردن شباهت بین کارها و نگاشت بین آنها که در دیگر روشهای انتقال دانش وجود دارد در اینجا وجود ندارد. همچنین روش پیشنهادی باعث بهبود پاداش بدست آمده و افزایش سرعت یادگیری می شود.

## مراجع

[1] Matthew E. Taylor and Peter Stone. Behavior Transfer for value function Based Reinforcement Learning. In conference on Autonomous agents and Multi agent System (AAMAS-05) pp. 53-59

[2] Yaxin Liu and Peter Stone. Value-Function- Based Transfer for Reinforcement Learning Using Structure Mapping. in Proceedings of the Twenty-First National Conference on Artificial Intelligence (AAAI-06),