

روشی نوین برای یادگیری تقویتی فازی باناظر برای ناوبری ربات

فاطمه فتحی نژاد^۱، ولی درهمی^۲

^۱ فارغ التحصیل کارشناسی ارشد مهندسی برق و کامپیوتر، گروه کامپیوتر، دانشگاه یزد fateme.fathinezhad@stu.yazduni.ac.ir

^۲ استادیار، دانشکده مهندسی برق و کامپیوتر، گروه کامپیوتر، دانشگاه یزد vderhami@yazduni.ac.ir

(تاریخ دریافت مقاله ۱۳۹۱/۴/۱۳، تاریخ پذیرش مقاله ۱۳۹۱/۷/۲)

چکیده: استفاده از یادگیری باناظر در ناوبری ربات‌های متحرک، با چالش‌های جدی از قبیل ناسازگاری و اختلال در داده‌ها، مشکل جمع‌آوری نمودن داده آموزش و خطای زیاد در داده‌های آموزشی مواجه می‌باشد. قابلیت‌های یادگیری تقویتی همچون عدم نیاز به داده آموزشی و آموزش تنها با استفاده از یک معیار اسکالر راندمان باعث کاربرد آن در ناوبری ربات شده است. از طرفی یادگیری تقویتی زمانبر بوده و دارای نرخ شکست‌های بالا در مرحله آموزش می‌باشد. در این مقاله، یک ایده جدید برای استفاده مؤثر از هر دو الگوریتم یادگیری فوق‌ارائه می‌شود. یک کنترلگر فازی سوگنو مرتبه صفر با تعدادی عمل‌کنندید برای هر قاعده جهت تولید فرمان‌های کنترل ربات در نظر گرفته شده است. هدف از آموزش تعیین عمل مناسب برای هر قاعده است. روش ترکیبی پیشنهاد شده دو مرحله دارد. در مرحله اول، داده آموزشی با حرکت ربات توسط ناظر در محیط جمع‌آوری می‌شود. سپس با بهره‌گیری از روش جدید ارائه شده، پارامترهای ارزش هر عمل‌کنندید در قواعد فازی با کمک داده‌های آموزشی مقداردهی اولیه می‌شوند. در مرحله دوم از الگوریتم سارسای فازی برای تنظیم دقیق‌تر پارامترهای تالی کنترلگر بصورت برخط استفاده می‌شود. نتایج شبیه‌سازی در شبیه‌ساز KIKS برای ربات خپرا حاکمی از بهبود قابل توجه در زمان یادگیری، تعداد شکست‌ها، و کیفیت حرکت ربات می‌باشد.

کلمات کلیدی: ناوبری ربات، یادگیری باناظر، یادگیری تقویتی، کنترلگر فازی.

A Novel Supervised Fuzzy Reinforcement Learning for Robot Navigation

Fateme Fathinezhad, Vali Derhami

Abstract: Applying supervised learning in robot navigation encounters serious challenges such as inconsistency and noisy data, difficulty to gathering training data, and high error in training data. Reinforcement Learning (RL) capabilities such as lack of need to training data, training using only a scalar evaluation of efficiency and high degree of exploration have encouraged researcher to use it in robot navigation problem. However, RL algorithms are time consuming also have high failure rate in the training phase. Here, a novel idea for utilizing advantages of both above supervised and reinforcement learning algorithms is proposed. A zero order Takagi-Sugeno (T-S) fuzzy controller with some candidate actions for each rule is considered as robot controller. The aim of training is to find appropriate action for each rule. This structure is compatible with Fuzzy Sarsa Learning (FSL) which is used as a continuous RL algorithm. In the first step, the robot is moved in the environment by a supervisor and the training data is gathered. As a hard tuning, the training data is used for initializing the value of each candidate action in the fuzzy rules. Afterwards, FSL fine-tunes the parameters of conclusion parts of the fuzzy controller online. The simulation results in KIKS simulator show that the proposed approach significantly improves the learning time, the number of failures, and the quality of the robot motion.

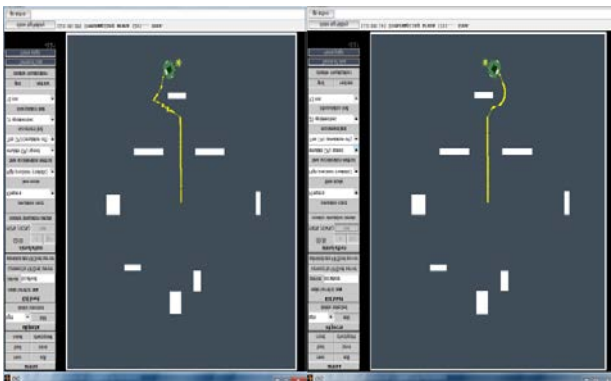
Keywords: Robot navigation, Supervised learning, Reinforcement learning, Fuzzy controller.

۱- مقدمه

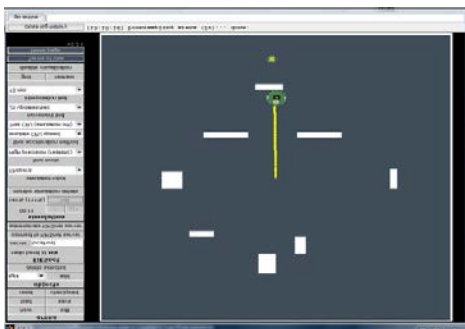
ناوبری برای ربات‌های متحرک عبارت است از حرکت از یک نقطه معین و رسیدن به یک هدف مشخص در حالی که ربات بتواند از برخورد به موانع اجتناب کند [۲،۱]. بطور کلی در یک محیط پویا، استفاده از الگوریتم‌های سراسری^۱ برای یافتن مسیر حرکت ربات، غیرممکن و یا بسیار پرهزینه است، زیرا در این روش‌ها مدل ریاضی یا نقشه کاملی از محیط مورد نیاز است. لذا چنانچه مشخصات محیط ناشناخته باشد و یا محیط در حال تغییر باشد، مسیریابی با استفاده از روش‌های طراحی مسیر محلی انجام می‌شود. روش‌های طراحی مسیر محلی از اطلاعات فراهم شده توسط حسگرهایی مانند حسگرهای سونار و یا حسگرهای مادون قرمز بهره می‌برند [۴،۳]. در میان روش‌های ارائه شده الگوریتم‌های هندسی فرض می‌کنند که حسگرهای سنجش فاصله نصب شده بر روی ربات قادرند، بطور کامل موانع را بصورت برخط تشخیص دهند. این دو فرض در محیط‌های واقعی غیرقابل قبول و برآوردن آنها وقت گیر است [۱]. از دیگر روش‌ها، روش‌های مبتنی بر پتانسیل می‌باشند که مؤثرتر از الگوریتم‌های هندسی به نظر می‌رسند، چرا که به جزئیات کمتری از موقعیت موانع نیاز دارند، لیکن این روش‌ها نیز نقاط ضعف زیر را دارند [۵،۱]:

- ۱- رخداد کمینه محلی منجر می‌گردد که ربات در حلقه ایجاد شده بین موانع به دام بیفتد.
 - ۲- حرکت ناپایدار ربات در کاربردهای عملی.
 - ۳- مشکل پیدا کردن ضرایب مؤثر مربوط به سرعت و نیرو در محیط‌هایی با موانع درهم که ارائه مدل ریاضی برای محیط را دشوار می‌کند.
 - ۴- افت راندمان به خاطر نایقینی و خطای مربوط به حسگرها (که برای هر حسگر متفاوت از بقیه است) به علت عملکرد بر اساس یک استراتژی از پیش تعیین شده و عدم وجود یادگیری.
- مشکلات اشاره شده در بالا، محققین را تشویق به استفاده از الگوریتم‌های یادگیرنده با استفاده از اطلاعات محلی در مسأله ناوبری ربات نموده است. این یادگیری با استفاده از اطلاعات حاصله از حسگرهای قرار گرفته بر روی ربات انجام می‌پذیرد [۶].
- یادگیری باناظر یکی از روش‌های قدیمی استفاده شده برای تنظیم پارامترهای کنترلگر می‌باشد که برای آموزش کنترلگر ربات نیز از آن استفاده شده است. در این روش ابتدا ربات در محیط توسط ناظر حرکت داده شده و سپس با توجه به داده‌های جمع‌آوری شده، با استفاده از روشهای مبتنی بر گرادینان [۷] پارامترهای کنترلگر در جهت کاهش مجموع مربعات خطای خروجی تنظیم می‌شوند.
- این الگوریتم در مسأله ناوبری ربات با ضعف‌های جدی مانند موارد زیر مواجه است:

- ۱- عدم اطلاع ناظر از فرمان کنترلی مناسب در بعضی از وضعیت‌ها: لذا در وضعیت‌های خاص خروجی تعیین شده توسط ناظر دارای خطای زیاد می‌باشد.
- ۲- ناسازگاری در داده‌ها: بعنوان مثال در نظر بگیرید یک ربات در جلوی مانع قرار دارد و سمت چپ و راست آن خالی می‌باشد. ناظر ممکن است در چنین وضعیتی یکبار با اعمال فرمان چرخش به راست (مثلا +۴۵ درجه) و یک بار دیگر با چرخش به چپ (مثلا -۴۵ درجه) ربات را از کنار مانع عبور دهد. این ناسازگاری باعث ایجاد مشکل در آموزش می‌شود. الگوریتم یادگیری باناظر برای این وضعیت یک عمل را باید تعیین نماید. واضح است که هر کدام را انتخاب کند خطا برای داده دیگر زیاد می‌شود و از آنجا که روش بر اساس کمینه کردن مجموع مربعات خطا می‌باشد، خروجی تعیین شده برای این وضعیت عددی نزدیک به صفر است. این خروجی به معنی حرکت مستقیم ربات به طرف جلو در این وضعیت و در نتیجه برخورد با مانع است. شکل ۱ بیانگر مورد مذکور می‌باشد. در دو قسمت "الف" و "ب" ربات توسط ناظر یکبار از سمت راست و یکبار از سمت چپ عبور داده شده است. نتیجه یادگیری با این داده ناسازگار حرکت مستقیم به سمت هدف و برخورد به مانع می‌باشد که در قسمت "ج" از شکل ۱ نشان داده شده است.



الف: چرخش به راست ب: چرخش به چپ



ج: حرکت مستقیم به طرف جلو و برخورد با مانع

شکل ۱: تاثیر سوء مشکل ناسازگاری در داده‌های آموزشی.

^۱-Global algorithms

ایده های ارائه شده در مراجع فوق برای فضای حالت گسسته و یا فضای عمل گسسته هستند، در حالیکه توجه ما در این مقاله بر روی فضای حالت - عمل پیوسته است.

در [۱] ایده ای برای فضای حالت و عمل پیوسته آمده است؛ در آن از داده آموزشی تولید شده توسط ناظر برای تنظیم اولیه پارامترهای بخش عملگر در معماری عملگر-نقاد^۳ [۱۵] استفاده شده است. در این مرجع ابتدا توسط یادگیری باناظر مقدار عمل برای هر حالت پیشنهاد می شود و سپس با استفاده از یادگیری تقویتی، مقدار نهای پیرامون مقدار پیشنهادی تنظیم می شود. روش مذکور دو ضعف عمده دارد:

۱- اثر سوء ناسازگاری داده اشاره شده در بالا، باعث خطای زیاد در خروجی تنظیم شده توسط یادگیری باناظر می گردد.

۲- ضعف عدم کاوش مناسب در معماری عملگر-نقاد [۱۵]، علاوه بر ضعف ذاتی معماری عملگر-نقاد در این خصوص از آنجا که مرحله تنظیم با یادگیری تقویتی تنظیمات پیرامون مقدار تنظیم شده با روش یادگیری باناظر صورت می گیرد، این ضعف تشدید شده است.

در کار قبلی، ما در مقاله [۱۶] ایده ای شبیه به روش فوق را در معماری نقاد-تنها بکار گرفتیم. بدین صورت که از داده های آموزشی با ناظر برای تنظیم اولیه توابع عضویت ورودی سیستم فازی استفاده شده است و آنگاه پارامترهای تالی کنترلگر فازی بصورت برخط با استفاده از روش یادگیری سارسای فازی^۴ (FSL) [۹] که یک روش یادگیری تقویتی فازی^۵ (FRL) با معماری نقاد-تنها^۶ است تنظیم شده است. هر چند مشکل عدم کاوش بخاطر استفاده از معماری نقاد-تنها مرتفع شده است لیکن هیچ آموزشی در خصوص عمل خروجی مناسب برای هر وضعیت صورت نگرفته است و در واقع مقادیر تالی کنترلگر فازی تنها با روش FSL تنظیم می گردند. بهمین دلیل بهبود بدست آمده فاحش نیست.

در اینجا روشی جدید برای ترکیب یادگیری باناظر و یادگیری تقویتی فازی با معماری نقاد-تنها ارائه می شود. لازم به ذکر است دو معماری معروف استفاده شده در FRL، معماری نقاد-تنها و عملگر-نقاد می باشند. از مزایای معماری نقاد-تنها پتانسیل بالا در برقراری تعادل بین کاوش و بهره برداری از تجربیات است. لذا این معماری برای مسائلی که نیاز به کاوش بالا دارند مانند ناوبری ربات مناسب می باشد. دو الگوریتم یادگیری سارسای فازی [۹] و یادگیری کیو فازی^۷ (FQL) [۴] بر اساس معماری نقاد-تنها ارائه شده اند. برای روش FQL نه تنها هیچ قضیه یا لمی در جهت همگرایی آن وجود ندارد بلکه مثالهای واگرایی [۹] آن نیز موجود می باشد. لیکن در [۹] قضایای مربوط به همگرایی و اثبات نقاط ایستای روش FSL آمده است لذا الگوریتم یادگیری تقویتی پیوسته استفاده شده در این مقاله روش FSL می باشد.

با توجه به ضعف های اشاره شده در یادگیری باناظر، استفاده از روش های هوشمند برای یادگیری ربات ها گسترش یافت. یادگیری تقویتی^۱ یک الگوریتم مدرن هوشمند است که به جهت دارا بودن قابلیت هایی همچون عدم نیاز به خروجی مطلوب، آموزش تنها با استفاده از یک معیار اسکالر راندمان، امکان آموزش برخط، و درجه کاوش بالا گزینه مناسبی جهت تنظیم پارامترهای کنترلگر ربات می باشد. در واقع در یادگیری تقویتی به عامل گفته نمی شود که عمل صحیح در هر وضعیت چیست، و فقط با استفاده از یک معیار اسکالر که سیگنال تقویتی نامیده می شود خوب یا بد بودن عمل به عامل نشان داده می شود. عامل موظف است با در دست داشتن این اطلاعات، یاد بگیرد که بهترین عمل کدام است. این ویژگی یکی از نقطه قوت های خاص الگوریتم یادگیری تقویتی است [۸]. اما از جنبه دیگر، دو چالش پیش روی یادگیری تقویتی زمانبر بودن و کند بودن آموزش در آن است. این مشکل در مسائل ناوبری ربات هم که معمولاً فضای حالت بزرگ است بطور جدی مشهود است. مسأله دیگر این است که امکان تنظیم همه پارامترهای کنترلگر (پارامترهای توابع عضویت مقدم در کنترلگرهای فازی یا پارامترهای وزن در لایه های ابتدایی کنترلگرهای عصبی) در الگوریتم های یادگیری تقویتی پیوسته که از آنها برای کنترلگر ربات استفاده شده است، وجود ندارد.

یک ایده سودمند برای بهره گیری از مزایا و کاهش ضعف های دو روش یادگیری تقویتی و یادگیری باناظر، استفاده از ترکیب این دو روش یادگیری می باشد. در [۱۱] از یادگیری باناظر برای تخمین اولیه احتمال انتخاب عمل استفاده شده است. نویسنده به دنبال روش ترکیبی از یادگیری تقویتی و یادگیری باناظر خطی است که از یادگیری باناظر خطی برای تولید سیاست انتخاب عمل در یادگیری تقویتی استفاده شده است. لذا انتخاب عمل در روش یادگیری کیو^۲ با توجه به احتمال انتخاب عمل هایی که از داده آموزشی بدست آمده است انجام می شود. روش مذکور در مسأله سیستم مکالمه بکار گرفته شده است. در [۱۲] یادگیری تقویتی باناظر برای مسئله دنبال کردن خط در ربات متحرک استفاده شده است و از دانش ناظر بعنوان دانشی که می تواند برای تصمیم در خصوص کاوش در مرحله انتخاب عمل استفاده شود، بهره برده شده است. در [۱۳] ناظر با استفاده از کنترلگر PID عملی را برای هر حالت انتخاب می نماید. سپس در هنگام انتخاب عمل، در روش یادگیری تقویتی عمل انتخاب شده توسط کنترلگر PID شانس بالاتری برای انتخاب خواهد داشت. تابع ارزش عمل هم بر اساس روش یادگیری کیو به روز رسانی می شود. در [۱۴] نیز از ترکیب یادگیری باناظر و یادگیری تقویتی برای مسأله حرکت ربات انسان نما به سمت شارژر و اتصال به آن استفاده شده است. در اینجا سعی شده ارزش اولیه عملها از یادگیری تقویتی استفاده شود. البته فضای عملها گسسته است و کلا چهار عمل برای ربات در نظر گرفته شده است.

³-Actor-Critic

⁴-Fuzzy Sarsa Learning

⁵-Fuzzy Reinforcement Learning

⁶-Critic-only

⁷-Fuzzy Q- Learning

¹-Reinforcement Learning

²-Q-learning

ارزش - عمل تقریب زده شده برای عمل a در حالت s که با $\tilde{Q}(s, a)$ نشان داده می شود، بصورت ذیل محاسبه می شوند [۹]:

$$a_t(s_t) = \sum_{i=1}^R \mu_i(s_t) o_{ii} \quad (1)$$

$$\tilde{Q}_t(s_t, a_t) = \sum_{i=1}^R \mu_i(s_t) w_i^{ij} \quad (2)$$

پس از محاسبه عمل نهایی a_t و اعمال آن، محیط به حالت جدید s_{t+1} رفته و عمل جدید a_{t+1} با توجه به مقادیر وزن فعلی w_t انتخاب می شود. ضمناً سیگنال تقویتی r_{t+1} از محیط دریافت می گردد. آنگاه مقادیر پارامترهای وزن هر قاعده بصورت زیر به روز رسانی می شوند [۹]:

$$\Delta w_{t+1}^{ij} = \begin{cases} \alpha_t \times \Delta \tilde{Q}_t(s_t, a_t) \times \mu_i(s_t) & \text{if } j = i^+ \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

که α نرخ آموزش و γ فاکتور نزول و $\Delta \tilde{Q}$ خطای تقاضای موقتی ارزش - عمل است و بصورت ذیل محاسبه می گردد:

$$\Delta \hat{Q}_t(s_t, a_t) = r_{t+1} + \gamma \hat{Q}_t(s_{t+1}, a_{t+1}) - \hat{Q}_t(s_t, a_t) \quad (4)$$

قدم های الگوریتم FSL بصورت خلاصه در زیر آمده است [۹]:

- ۱- مشاهده حالت s_{t+1} و دریافت سیگنال تقویتی r_{t+1} .
- ۲- انتخاب یک عمل مناسب در هر قاعده با استفاده از روش انتخاب عمل بیشینه نرم.
- ۳- محاسبه عمل نهایی a_{t+1} و مقدار تقریبی تابع ارزش - عمل $\hat{Q}_t(s_{t+1}, a_{t+1})$ با استفاده از (۱) و (۲).
- ۴- محاسبه $\Delta \hat{Q}$ و بروزرسانی w با استفاده از (۳) و (۴).
- ۵- محاسبه مقدار تقریبی جدید $\hat{Q}_{t+1}(s_{t+1}, a_{t+1})$ با استفاده از (۴)
- ۶- اجرای عمل نهایی.
- ۷- $t \leftarrow t + 1$ و بازگشت به مرحله اول.

۳- طراحی کنترلر برای ناوبری ربات

برای بررسی ایده مورد بحث از ربات خپرا^۱ که یک ربات مینیا توری برای فعالیت های آزمایشگاهی و تحقیقاتی ساخته شده توسط شرکت سویسی K-Team [۱۷] است، استفاده می کنیم. قابلیت فراوان این ربات و اندازه مناسب آن جهت فعالیت های آزمایشگاهی، منجر به استقبال گسترده محققین در استفاده از این ربات در ارزیابی روش های خود شده است [۱۸].

پیرامون ربات خپرا هشت حسگر مادون قرمز که هر یک دارای یک فرستنده و گیرنده هستند، وجود دارد. هر دو حسگر در یک وجه ربات نصب شده است (شکل ۲). محدوده عملکرد مؤثر حسگرهای این ربات

ایده آن است که با کمک دانش ناظر بجای تعیین یک عمل برای هر حالت، ارزش اولیه برای اعمال ممکن کنترلر تعیین می شود. سپس با کمک یادگیری تقویتی، بصورت برخط تنظیم نهایی درجهت بهبود کارایی صورت می گیرد. این ترکیب باعث ایجاد تسریع در فرایند یادگیری، بهبود کیفیت آموزش، و کاهش تعداد برخوردهای ربات به موانع و همگرایی سریعتر در حین آموزش می شود. براساس بررسی های ما، این کار اولین روش ارائه شده برای ترکیب یادگیری باناظر و یادگیری تقویتی فازی پیوسته با معماری نقاد- تنها می باشد.

ساختار مقاله به شرح زیر است. در بخش دوم الگوریتم FSL شرح داده می شود. بخش سوم نحوه طراحی ساختار کنترلر فازی را شرح می دهد. در بخش چهارم ایده مقاله برای ترکیب یادگیری تقویتی و یادگیری باناظر را شرح می دهیم. بخش پنجم به پیاده سازی و شبیه سازی کار پرداخته است. در بخش آخر بحث و نتیجه گیری آمده است.

۲- یادگیری سارسای فازی (FSL)

از آنجا که FSL از لحاظ تحلیل ریاضی و عملکرد نسبت به FQL ارجحیت دارد، بعنوان الگوریتم پایه در فرآیند یادگیری در کار ما انتخاب شده است. این الگوریتم برخلاف FQL که مستقل از سیاست می باشد، یک روش وابسته به سیاست است که تالی قاعده سیستم فازی را بصورت برخط تنظیم می نماید.

الگوریتم FSL از ترکیب سیستم های فازی بعنوان یک تقریب زننده تابعی خطی با روش سارسا [۱۰] حاصل شده است. یک سیستم فازی سوگنو مرتبه صفر با n ورودی و یک خروجی و R قاعده به فرم زیر را در نظر بگیرید:

Ri : If x_1 is L_{i1} and \dots and x_n is L_{in} then

$$\begin{array}{ll} o_{i1} & \text{with value } w^{i1} \\ \text{or } o_{i2} & \text{with value } w^{i2} \\ \vdots & \vdots \\ \text{or } o_{im} & \text{with value } w^{im} \end{array}$$

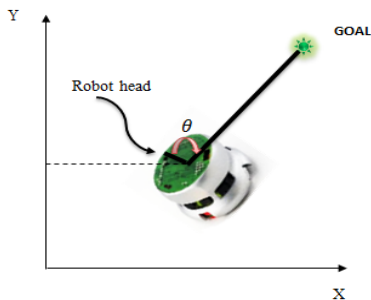
که در آن $s = x_1 \times \dots \times x_n$ بردار n بعدی متغیرهای حالت ورودی، $L_i = L_{i1} \times \dots \times L_{in}$ شامل n مجموعه فازی محدب نرمال با مرکزهای یکتا برای i امین قاعده، m تعداد عمل های گسسته ممکن برای هر قاعده، o_{ij} i امین عمل کاندید در قاعده i ام و w^{ij} مقدار ارزش تقریب زده شده برای عمل i ام در قاعده i ام است. در هر قدم زمانی برای هر قاعده، یک عمل از میان m عمل کاندید شده برای تالی قاعده بر مبنای مقدار وزن آن عمل انتخاب می شود و آنگاه عمل نهایی از ترکیب وزن دار این عمل ها حاصل می گردد. هدف آموزش به روز رسانی بر خط مقادیر وزن w^{ij} با توجه به سیگنال تقویت دریافت شده است، بگونه ای است که بهترین انتخاب عمل بر مبنای آن ها حاصل گردد [۹].

شدت آتش هر قاعده از حاصلضرب درجه های تطابق مقدم قاعده برای ورودی های مختلف بدست می آید و خروجی سیستم a و مقدار

^۱-Khepera robot

به اندازه زوایه مذکور چرخیده و پس از همراستا شدن بطور مستقیم به سمت هدف می رود. همان طور که در شکل ۳ نیز مشخص می باشد، زمانی که ربات نزدیک مانع می شود رفتار "پیگیری هدف" غیرفعال می شود و خروجی حاصل از ماژول "اجتناب از موانع" به ربات اعمال می شود. وظیفه این ماژول تعیین زاویه حرکت ربات در هر قدم زمانی به گونه ای است که ضمن پرهیز از برخورد به موانع در جهت نزدیک شدن به هدف، ربات حرکت کند. توجه شود از آنجا که در این ماژول خروجی تولید شده با در نظر گرفتن دو مورد اجتناب از موانع و نزدیک شدن به هدف تولید می شود. لذا دیگر مانند دیگر کارهای مرتبط در این زمینه [۲۰] نیازی به یک ماژول برای ترکیب خروجی های رفتارها نیست.

از این روزبه محاسبات و پیچیدگی سیستم کاهش یافته است. پیشنهاد ما برای طراحی ساختار این کنترلگر، یک کنترلگر فازی سوگنو مرتبه صفر می باشد. ساختار این کنترلگر بصورتی در نظر گرفته می شود که با ساختار استفاده شده در FSL همخوانی داشته باشد.



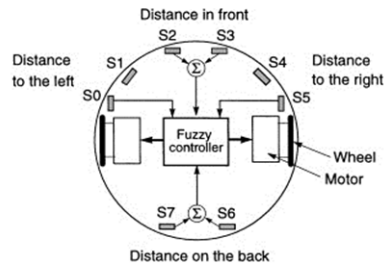
شکل ۴: محاسبه زاویه چرخش ربات برای رفتار پیگیری هدف.

کنترلگر مورد نظر دارای چهار ورودی (سه ورودی اول بعنوان معیار فاصله ربات با مانع در یکی از سه جهت راست، جلو، و عقب هستند، و ورودی چهارم زاویه پیشانی ربات با هدف) و یک خروجی (مقدار زاویه چرخشی پیشانی ربات در هنگام نزدیکی به موانع) می باشد. توابع عضویت ورودی این کنترلگر بشکل گوسی در نظر گرفته شده اند. در هر بعد ورودی بترتیب ۲، ۳، و ۲ مجموعه فازی تعریف شده است. با توجه به تقسیم بندی انجام شده در هر بعد ورودی، کنترلگر دارای ۲۴ قاعده می باشد. مقدار خروجی هر قاعده یک مقدار ثابت است که باید از مجموعه عملهای کاندید در نظر گرفته شده برای هر قاعده $A = \{O_{i1}, O_{i2}, \dots, O_{im}\}$ انتخاب شود. عمل مناسب برای هر قاعده از این مجموعه عمل کاندید تعیین می گردد. لذا هدف از آموزش تعیین عمل مناسب از میان مجموعه عمل کاندید برای تالی هر قاعده است.

۴- یادگیری سارسای فازی باناظر

در این بخش روش جدیدی برای تعیین عمل مناسب از میان مجموعه عملهای کاندید ممکن برای تالی هر قاعده، در ساختار کنترلگر فازی سوگنو مرتبه صفر ارائه می گردد. روش ارائه شده که ترکیبی از یادگیری باناظر و یادگیری تقویتی است شامل دو مرحله می باشد:

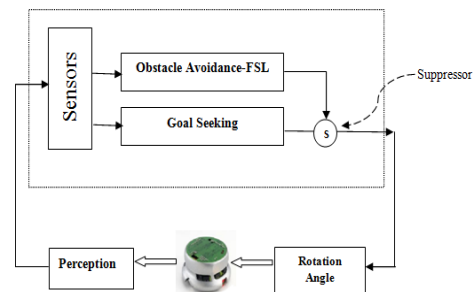
بین ۱ تا ۵ سانتیمتر است و مقدار خروجی آنها با فاصله ربات تا مانع رابطه عکس دارد. هر چه ربات به مانع نزدیک تر گردد، مقدار خروجی حسگر بیشتر و هر چه دورتر گردد، مقدار کمتری خواهد داشت.



شکل ۲: موقعیت سنسورها در ربات میناتور خپرا.

بر روی هر چرخ ربات یک رمزگذار نصب شده است که با شمارش پالس های حاصل از رمزگذارها می توان مسافت طی شده توسط هر چرخ را محاسبه نمود. همچنین بسته به کاربرد می توان تجهیزات جانبی دیگری همچون دوربین، چنگک و غیره بر روی ربات نصب نمود [۱۷].

در اینجا هدف آن است که اگر ربات در مجاورت موانع قرار دارد، بدون برخورد به مانع با توجه به موقعیت هدف از موانع عبور کند، و زمانی که پیرامون ربات مانعی وجود ندارد، ربات بطور حریصانه، به سمت هدف چرخیده به سمت آن حرکت کند. برای عملکرد بهتر ربات هنگام حرکت و نیز کاهش پیچیدگی سیستم، از معماری رده بندی ۲ که بروکس در [۱۹] ارائه نموده است، استفاده می نمایم. با استفاده از این معماری برای ربات دو رفتار در نظر گرفته می شود. یک رفتار بنام "اجتناب از موانع" برای زمانهایی که ربات نزدیک موانع است و رفتار دیگر "پیگیری هدف" برای زمانی که پیرامون ربات مانعی وجود ندارد.



شکل ۳: معماری مبتنی بر رفتار.

شکل ۳ طراحی انجام شده در این مقاله بر مبنای معماری مذکور را نشان میدهد. خروجی رفتار "پیگیری هدف" همان طور که در شکل ۴ آمده است به سادگی با محاسبه زاویه پیشانی ربات با هدف بدست می آید. بدین معنی که ابتدا مقدار اختلاف زاویه بدست آمده و سپس ربات

۱-Encoder

۲-Subsumption

- ۷- مراحل بالا را بطور کامل برای μ_{12} و سپس برای μ_{13} (مانند قبل با جایگزینی مقدار μ_{13} بجای μ_{12}) و μ_{14} (با جایگزینی مقدار μ_{14} بجای μ_{13}) تکرار می کنیم.
- ۸- در نهایت وقتی که مراحل بالا برای همه نمونه داده های جمع آوری شده انجام شد، ارزش عمل j امین از قاعده i امین بصورت زیر مقداردهی اولیه می کنیم.

$$w^{ij} = \frac{c_{ij}}{(\sum_j c_{ij})^2} \quad (6)$$

شبه کد روش پیشنهاد شده در شکل ۵ آمده است.

```

Initialize counter ( $C_{ij}$ ) for  $j$ -th candidate action in  $i$ -th rule with zero, assign set  $A_{in}$ 
Foreach ( $x_p, y_p$ ),  $x_p$  is input of the controller and  $y_p$  is suggested output by supervisor do
  Select four rules that have highest firing strength ( $l_1, l_2, l_3, l_4$ )
  output =  $y_p, i=1$ 
  repeat for each rule  $L_i$ 
    Calculate firing strength ( $\mu_i$ )
     $a = \lfloor \frac{output}{\mu_i} \rfloor$ 
    closest action to the division result is selected and indexed by  $k_i$ 
     $c_{ij} = c_{ij} + 1$ 
     $y'_p = output - \mu_i \times a_{ilk1}$ 
    output =  $y'_p$ 
  until  $i=4$ 
end
foreach rule  $i$ 
  foreach candidate action  $j$ 
     $w^{ij} = \frac{c_{ij}}{(\sum_j c_{ij})^2}$ 
  end
end

```

شکل ۵: شبه کد روش پیشنهاد شده برای یافتن ارزش اولیه عمل های کاندید.

پس از تعیین مقدار اولیه w^{ij} ها در مرحله دوم از الگوریتم FSL برای تنظیم برخط تالی قواعد کنترلر فازی که مقدار ارزش عمل های آن (w^{ij}) بصورت بالا مقداردهی اولیه شده است، استفاده می کنیم. روش ترکیبی مذکور را یادگیری سارسای فازی باناظر^۱ (SFSL) می نامیم. بلوک دیاگرام SFSL در شکل ۶ آمده است. بطور خلاصه روش SFSL شامل مراحل زیر می شود:

- ۱- حرکت ربات در محیط و جمع آوری داده های آموزشی.
- ۲- مقداردهی اولیه ارزش عمل های کاندید با روش ارائه شده (شکل ۵).
- ۳- تنظیم نهایی مقدار تالی قواعد با استفاده از FSL.

در مرحله اول ابتدا با حرکت ربات در محیط توسط ناظر، داده آموزشی جمع آوری می شود. در اینجا بر خلاف روش های موجود که از داده آموزشی برای تعیین عمل مشخص برای هر حالت استفاده می شود، یک روش جدید جهت استفاده از داده آموزش برای ارزش گذاری عمل های ممکن در هر حالت ارائه می شود. بدین صورت که این داده آموزشی برای مقداردهی اولیه ارزش هر عمل کاندید w^{ij} (معرفی شده در بخش دوم) در تالی هر قاعده کنترلر فازی استفاده می شود. از این رو هدفی که ما بدنبال آن هستیم تعیین ارزش برای هر خروجی انتخاب شده در هر حالت توسط ناظر است.

بدین مفهوم که مثلاً اگر ناظر عملهای متفاوتی را در یک وضعیت خاص در دفعات مجزا انتخاب کند، متناسب با تعداد انتخاب هر عمل در آن وضعیت خاص، به آن عمل ارزش داده شود. از آنجا که خروجی نهایی سیستم فازی از ترکیب وزن دار تالی انتخاب شده در هر قاعده بدست می آید. لازم است برای هر خروجی یک ترکیب ممکن از عمل های کاندید هر قاعده به گونه ای که ترکیب این اعمال بتواند منجر به مقداری نزدیک به آن خروجی شود، پیدا نموده و آنگاه ارزش آن عمل ها افزایش یابد.

هر نمونه p ام از داده های جمع آوری شده شامل جفت داده ورودی- خروجی (x_p, y_p) را در نظر بگیرید که x_p ورودی کنترلر و y_p خروجی پیشنهاد شده توسط ناظر می باشد. قدمهای زیر برای تعیین ارزش اولیه عمل های کاندید در هر قاعده (w^{ij}) ، برای هر نمونه p ام از داده ها (x_p, y_p) دنبال می شود.

- ۱- برای ورودی x_p ، چهار قاعده غالب (قاعده هایی که بیشترین میزان شدت آتش (μ) را دارند) را انتخاب می نماییم. این قواعد با سمبلهای l_1 و l_2 و l_3 و l_4 نشان داده می شوند، طوری که:

$$\mu_{11} < \mu_{12} < \mu_{13} < \mu_{14}$$

- ۲- y_p بر μ_{11} (بیشترین شدت آتش) تقسیم می شود.

- ۳- نتیجه تقسیم با هر یک از عمل های کاندید (O_{11j}) مقایسه می شود. سپس نزدیک ترین عمل به نتیجه تقسیم انتخاب شده و k_1 بعنوان اندیس آن عمل در نظر گرفته می شود.

- ۴- شمارنده C_{ij} را برای نشان دادن دفعات انتخاب j امین عمل کاندید، در i امین قاعده بکار می بریم. در این مرحله مقدار شمارنده عمل k_1 ام در قاعده l_1 ، (C_{11k1}) یک واحد افزایش می یابد.

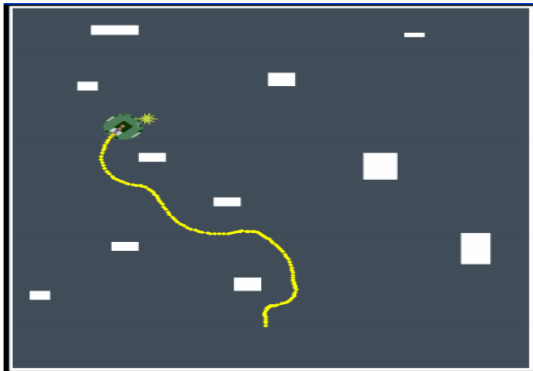
- ۵- عمل k_1 ام از مجموعه عمل های کاندید قاعده l_1 در شدت آتش این قاعده (μ_{11}) ضرب شده و این حاصل ضرب را از y_p کم می کنیم.

$$(y'_p = y_p - \mu_{11} \times a_{11k1}) \quad (5)$$

- ۶- مقدار y'_p را با مقدار y_p اولیه و مقدار شدت آتش قاعده l_2 (μ_{12}) را با مقدار شدت آتش قاعده l_1 (μ_{11}) جایگزین

می کنیم.

¹-Supervised Fuzzy Sarsa Learning



شکل ۷: نمونه ای از محیط آموزش در شبیه ساز KIKS.

با توجه به مطالب ذکر شده در بخش سوم، معماری رده‌بندی برای حرکت ربات در این مقاله پیشنهاد شد و هدف اصلی آموزش، تنظیم پارامترهای کنترلگر فازی برای مازول "اجتناب از موانع" در معماری طراحی شده شکل ۳، می‌باشد. سه ورودی اول تعریف شده در بخش سوم برای کنترلگر ربات با ترکیب خروجی های حسگرهای مادون قرمز هر کدام از وجوه راست، جلو و چپ بصورت زیر حاصل می‌شوند:

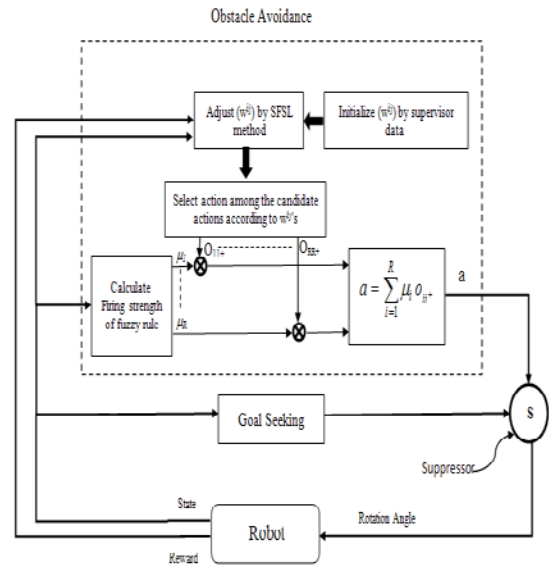
$$d_{face}(t) = 1 - \frac{\max(s_{face,1}, s_{face,2})}{1024} \quad (7)$$

$face \in \{Left, Front, Right\}$

که $s_{face,1}$ و $s_{face,2}$ مقدار خروجی حسگرهای یک و دو برای وجه مشخص شده در اندیس آن (چپ، جلو، و یا راست) می‌باشند. لازم به ذکر است خروجی حسگرها که در شبیه ساز عددی بین ۰ تا ۱۰۲۳ است بصورت پیش فرض همراه با مقداری نویز جمع شده است. مقدار صفر برای وقتی است که حسگر هیچ مانعی را در محدوده حس نکند و مقدار ۱۰۲۳ برای هنگامی است که حسگر تقریباً به مانع چسبیده است. ورودی چهارم زاویه پیشانی ربات با هدف است، که آن را با $\theta(t)$ نشان می‌دهیم و مقدار آن عددی بین ۱۸۰ و -۱۸۰ درجه می‌باشد.

چهار ورودی مذکور نرمالیز شده و سپس به کنترلگر وارد می‌گردند. خروجی کنترلگر زاویه‌ی چرخش پیشانی ربات است که عددی بین -۴۵ درجه تا +۴۵ درجه در نظر گرفته شده است.

مرحله اول شامل حرکت ربات توسط ناظر در محیط آموزش شکل ۷ با کمک جوی استیک^۳ می‌باشد. ۱۲۰۰ جفت داده ورودی - خروجی در این مرحله جمع آوری شد. از این داده ها با استفاده از روش ارائه شده در بخش چهارم، مقدار اولیه ارزش عمل‌های کاندید ($w^{(j)}$) در هر قاعده کنترلگر فازی مشخص گردید.



شکل ۶: نمودار بلوکی روش SFSL.

۵- شبیه سازی

شبیه‌سازهای متعددی برای ربات خپرا موجود است. در این میان، شبیه‌ساز KIKS که یک شبیه‌ساز ربات خپرا در محیط برنامه نویسی MATLAB است، برای مسأله ناوبری ربات استفاده می‌شود که این شبیه‌ساز مورد توجه بسیاری از محققین قرار گرفته است [۲۱]. در این پژوهش نیز از این شبیه‌ساز استفاده شده است. در ابتدای کار لازم است محیط‌های لازم برای شبیه سازی آماده گردد. برای هر محیط ابعاد آن، موقعیت و شکل موانع، موقعیت شروع حرکت ربات، و موقعیت هدف تعیین می‌گردد. برای این کار از واسط گرافیکی شبیه‌ساز و دستورات مرتبط با آن در شبیه‌ساز استفاده شده است.

ابعاد محیط آموزشی 820×820 میلی متر مربع است که موانع مختلفی با شکل‌های متفاوت در آن قرار گرفته‌اند. شکل ۷ محیط آموزش ربات را با موقعیت‌های مختلف شروع حرکت ربات و هدف نشان می‌دهد. در این محیط، مستطیل‌های سفید رنگ موانع و ستاره زرد رنگ، هدف می‌باشد.

هر رویداد^۱ در بخش آموزش شامل شروع از مبدأ و حرکت ربات تا رسیدن به هدف است. موقعیت هدف و شروع حرکت ربات در هر رویداد متفاوت می‌باشد. بخش آموزش در صورتیکه ربات به کران بالای تعداد حرکت‌ها که ۵۰۰ است و یا اینکه ۱۰ بار بطور متوالی بدون شکست به هدف برسد، به پایان می‌رسد. شماره رویدادها در پایان بخش آموزش بعنوان معیار زمان آموزش^۲ LDI در نظر گرفته می‌شود.

^۱-Episode

^۲-Learning Duration Index

^۳-Joystick

جهت مقایسه روش ارائه شده با یک روش ترکیبی مشابه، الگوریتم CSLAFSL [۱۶] انتخاب شد. برای روش CSLAFSL از داده آموزشی بدست آمده توسط حرکت ربات توسط ناظر برای تنظیم درجه عضویت توابع عضویت ورودی قواعد فازی طبق روش بیان شده در مقاله [۱] استفاده شد و آنگاه در معماری رده بندی ارائه شده (شرح داده شده در بخش سوم) بکار گرفته شد. همچنین در روش FSL [۹]، و FQL [۴] نیز در معماری رده بندی ارائه شده بکار رفتند و نتایج شبیه سازی آورده شده است. توجه شود که در این دو روش ارزش اولیه عمل‌های تالی ها ($w^{(j)}$) صفر می باشند (مقدار دهی اولیه نشده‌اند).

نتایج شبیه سازی در جدول ۱ آورده شده است. ستون اول این جدول چهار الگوریتم یادگیری را نشان می‌دهد. ستون دوم این جدول متوسط LDI ها را در بخش آموزش نشان می‌دهد که مقدار آن از متوسط‌گیری بر روی ۱۰ اجرای مستقل بدست آمده است. ستون سوم و چهارم نشان دهنده متوسط تعداد برخوردهای ربات با موانع، بترتیب در بخش آموزش و تست می‌باشد. نهایتاً در ستون پنجم متوسط مسافت پیموده شده توسط ربات در بخش تست آورده شده است.

جدول ۱: نتایج شبیه سازی در مسأله ناوبری ربات.

Methods	Ave. LDI	Failure Rate 1	Failure Rate 2	Ave. Distance
SFSL	40	30.4	8.7	78.00
CSLAFSL	88	49.3	9.1	78.04
FSL	107	62.6	9.3	78.12
FQL	124	66.5	9.8	77.08

همانطور که از نتایج مشهود است عملکرد روش SFSL بطور قابل توجهی از سه روش دیگر بهتر است. این روش برای معیار Ave. LDI نشانگر سرعت آموزش است، ۵۰ درصد بهتر از CSLAFSL، ۶۲ درصد سریعتر از FSL و ۶۸ درصد سریعتر از FQL، کنترلگر فازی را تنظیم می‌کند. عبارتی این روش سرعت زمان آموزش را حداقل ۶۰ درصد افزایش داده است. همچنین تعداد شکست‌ها در بخش آموزش در روش SFSL از سه روش بطور قابل توجهی (تقریباً ۵۰ درصد) کمتر شده است. تعداد شکست‌ها در بخش تست نیز در روش SFSL، ۴ درصد کمتر از CSLAFSL، ۶ درصد کمتر از FSL و ۱۱ درصد کمتر از FQL می‌باشد. از آنجا که هر چهار روش در همه تکرارها در مرحله تست به هدف رسیده‌اند مسافت طی شده تا هدف در آنها تقریباً یکسان است. جهت نمایش نحوه تغییرات مقادیر وزن عمل‌های کاندید، نمودار تغییرات مقدار ارزش ($w^{(j)}$) عمل‌های کاندید در قاعده ۲۳ ام کنترلگر فازی در شکل ۹ آورده شده است. همانطور که دیده می‌شود ارزش مربوط به اولین عمل کاندید (مربوط به 45°) بیشترین مقدار را دارد که پس از گذشت زمان کوتاهی از آموزش مقدار ارزش آن از بقیه عمل‌ها پیشی گرفته است. همچنین

در مرحله دوم آموزش از الگوریتم FSL، برای تنظیم برخظ پارامترهای کنترلگر استفاده می‌شود. در الگوریتم FSL معیار فاصله ربات با موانع که d نامیده می‌شود، بصورت زیر تعریف گردید:

$$d = \min(d_{face}) \quad face \in \{Left, Front, Right\} \quad (8)$$

هرگاه d صفر شود یک شکست^۱ به حساب می‌آید. هرگاه فاصله مرکز ربات تا هدف به ۵۰ میلی متر برسد، به معنی رسیدن ربات به هدف است. سیگنال تقویتی را با توجه به نزدیکی به موانع و زاویه سر ربات با هدف بصورت زیر تعریف می‌نماییم:

$$r(t) = \begin{cases} -1 & failure \\ -0.5 & d < 0.075 \\ \Delta/150 & \Delta > 0 \text{ \& } d \geq 0.075 \\ -0.01 + \Delta/150 & \Delta \leq 0 \text{ \& } d \geq 0.075 \\ 1 & goal \end{cases} \quad (9)$$

$$\Delta = |\theta(t-1)| - |\theta(t)|$$

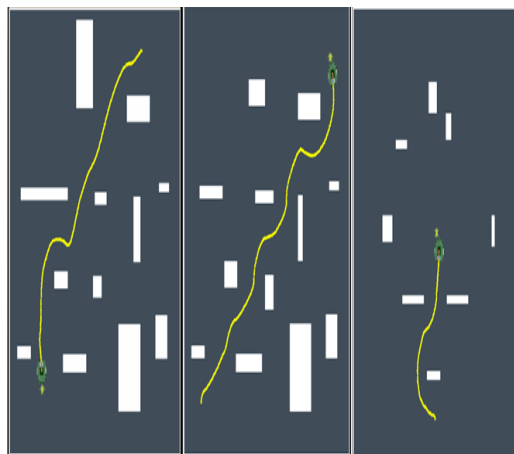
برای تالی هر قاعده ۱۳ عمل کاندید بصورت زیر در نظر گرفته شده-

است.

$$A = \{-45, -30, -20, -15, -10, -5, 5, 15, 10, 20, 30, 45\}$$

برای آموزش در این مرحله، پانصد جفت موقعیت تصادفی برای نقطه شروع حرکت ربات و هدف در محیط آموزشی شکل ۷ تولید شدند. ده اجرای مستقل انجام پذیرفت. هر اجرا از دو بخش آموزش و تست تشکیل می‌گردد.

پس از پایان آموزش، برای ارزیابی عملکرد ربات، بخش تست انجام می‌پذیرد. در این بخش ۱۰ محیط جدید ایجاد شده که در ۸ محیط اول تنها موقعیت شروع حرکت ربات و هدف متفاوت است. شکل ۸ این محیط‌ها را نشان می‌دهد. کیفیت عملکرد در محیط تست با معیارهای تعداد برخورد به موانع، و مسافت طی شده تا رسیدن به هدف ارزیابی می‌شود.



شکل ۸: محیط‌های تست مختلف برای رویدادهای ۱ تا ۱۰ در بخش تست

^۱- Failure

ماژول برای ترکیب رفتارها ندارد و در هر لحظه تنها خروجی یک رفتار به ربات اعمال می شود، لذا هزینه محاسباتی و طراحی کاهش یافته است. در خصوص تحلیل ریاضی روش ارائه شده باید گفته شود که از آنجا که روش مذکور در واقع از داده باناظر برای مقدار دهی اولیه پارامترهای روش FSL استفاده کرده است. لذا تمام شرایط بیان شده در قضایای ارائه شده برای FSL در مقاله [۹] را دارد و تحلیلها و قضایای ریاضی بیان شده در آن مرجع برای SFSL نیز برقرار است.

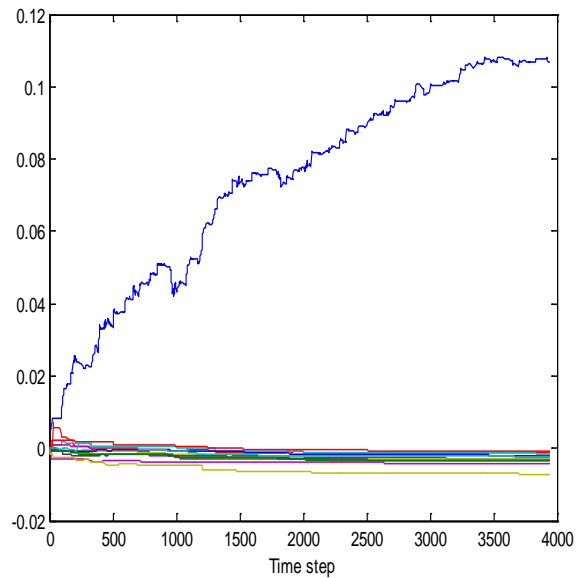
سپاسگزاری

این تحقیق با حمایت "صندوق حمایت از پژوهشگران کشور INSF" انجام شده است که بدینوسیله از آن مرکز محترم تشکر و قدردانی می گردد.

مراجع

- [1] C. Ye, N. H. C. Yung, and D. Wang, "A fuzzy controller with supervised learning assisted reinforcement learning algorithm for obstacle avoidance," *IEEE Transaction Systems, Man, Cybernetics*, vol. 33, no. 1, pp.17-27, Feb. 2003.
 - [2] T. Belker, M. Beetz, and A. Cremers, "Learning action models for the improved execution of navigation plans," *Robotics and Autonomous Systems*, vol. 38, pp. 137-148, Mar. 2002.
 - [3] T. Fong, I. Nourbakhsh, and K. Dautenhahn, "A survey of socially interactive robots," *Robotics and Autonomous Systems*, vol. 42, pp. 143-166, Mar. 2003.
 - [4] L. Jouffe, "Fuzzy inference system learning by reinforcement methods," *IEEE Trans. Syst., Man, Cybern. C*, vol. 28, no.3, pp. 338-355, Aug. 1998.
 - [5] H. R. Beom, and H. S. Cho, "A sensor-based navigation for a mobile robot using fuzzy logic and reinforcement learning," *IEEE Trans. Syst., Man, Cybern.*, vol. 25, no. 3, pp. 464-477, Mar. 1995.
 - [6] K. Macek, I. Petrovic, and N. Peric, "A reinforcement learning approach to obstacle avoidance of mobile robots," *Proc. IEEE Int. Conf. Advanced Motion Control*, vol.1, pp. 462-466, 2002.
 - [7] J. S. R. Jang, C. T. Sun, and E. Mizutani, "Neuro-Fuzzy and soft computing," Prentice-Hall, 1997.
- [۸] م. کلامی هریس، ن. پریش، م. ب. نقیبی سیستمی "بررسی یادگیری تقویتی و خواص سیاست بهینه در مسائل جدولی با استفاده از روش های کنترل دیجیتال"، مجله کنترل، جلد ۳، شماره ۱، بهار ۱۳۸۸
- [9] V. Derhami, V. Majd, and M. Nili Ahmadaabaadi "Fuzzy Sarsa learning and the proof of its

ترتیب مقدار ارزش عملهای دیگر نیز پس از گذشت زمان کوتاهی از آموزش دیگر تغییر نمی کند.



شکل ۹: نمودار تغییرات مقدار ارزش عملها در تالی قاعده ۲۳ ام کنترلگر

۶- بحث و نتیجه گیری

در این مقاله یک روش جدید برای ترکیب یادگیری باناظر و یادگیری تقویتی پیوسته پیشنهاد شد. جهت بررسی عملکرد روش مذکور از مسأله ناوبری ربات استفاده شد. در این روش فرمان های کنترلی هدایت ربات، براساس معماری رده بندی بروکس طراحی شد و فرمان نهایی از ترکیب دو رفتار "اجتناب از موانع" و "پیگیری هدف" بدست آمد. برای رفتار "اجتناب از موانع" یک کنترلگر فازی سوگنو مرتبه صفر طراحی شد. هدف یادگیری یافتن مقدار تالی مناسب برای هر قاعده این کنترلگر بود. روش ارائه شده از داده های آموزشی تولید شده که با کمک ناظر و از طریق حرکت دادن ربات در محیط بدست آمده بود، برای تقریب مقدار ارزش هر عمل کاندید استفاده نمود. در مرحله دوم آموزش، از الگوریتم FSL به عنوان الگوریتم یادگیری تقویتی پیوسته برای تنظیم نهایی مقدار تالی قواعد کنترلگر فازی بهره برده شد. نتایج شبیه سازی برای روش ارائه شده نشان داد که زمان آموزش و تعداد برخورد به موانع، کاهش قابل توجهی نسبت به سه روش CSLAFSL، FSL و FQL هنگامی که پارامترهای آنها مقدار دهی اولیه نشده اند، دارد. لازم به توجه است که در اینجا برخلاف روش های مرسوم که از یادگیری باناظر برای تعیین مقدار خروجی کنترلگر برای هر حالت استفاده می کنند، از داده های آموزشی بدست آمده، برای ارزش دهی به عمل های کاندید در هر حالت استفاده شد. بدین طریق نه تنها از اثر مخرب داده های ناسازگار جلوگیری بعمل آمد، بلکه از دانش موجود در این داده ها نیز سود برده شد. نکته قابل توجه دیگر استفاده از معماری رده بندی جهت ناوبری ربات بود. ساختار ارائه شده دیگر نیازی به یک

- [10] R. S. Sutton, and A. G. Barto, "Reinforcement learning: An introduction," Cambridge, MIT Press, 1998.
- [۱۶] ف. فتحی نژاد، و. درهمی "ترکیب یادگیری با ناظر با یادگیری تقویتی برای ناوبری ربات"، هفدهمین کنفرانس ملی سالانه انجمن کامپیوتر ایران، صفحه: ۱۱۵-۱۱۹، اسفند ۱۳۹۰
- [17] H. Maaref, and C. Barret "Sensor-based navigation of a mobile robot in an indoor environment" Robotics and Automation systems, vol. 38, pp. 1-18, Jan. 2002.
- [18] E. O. Ari, I. Erkmen, and A. M. Erkmen, "A FACL controller architecture for a grasping snake robot," Proc. IEEE Int. Conf. Intelligent Robots and Systems, pp. 1748-1753, 2005.
- [19] R. A. Brooks, "A robust layered control system for a mobile robot," Journal of Robotics and Automation, vol. 2, pp. 14-23, Mar. 1986.
- [20] K. Anam, Prihastono, H. Wicaksono, R. Effendi, S. Kuswadi, "Hybridization of Fuzzy Q-learning and Behavior-Based Control for Autonomous Mobile Robot Navigation in Cluttered Environment" International Joint ICROS- Conf SICE, pp. 1023 - 1028, Aug. 2009.
- [21] T. Nilsson, "KIKS: KIKS Is a Khepera Simulator" [http:// www.kiks.f2s.com](http://www.kiks.f2s.com)
- stationary points" Asian Journal of Control, vol. 10, No. 5, pp. 535-549, September 2008.
- [11] J. Henderson, O. Lemon, K. Georgila, "Hybrid Reinforcement/Supervised Learning for Dialogue Policies from communicator data," In IJCAI workshop on Knowledge and Reasoning in Practical Dialogue Systems, 2005.
- [12] R. Iglesias, C. V. Regueiro, J. Correa, S. Barro, "supervised reinforcement learning: application to a wall following behaviour in a mobile robot," Lecture Notes in Computer Science, vol. 1416, pp. 300-309, 1998.
- [13] L. Lin, H. Xie, D. Zhang, L. Shen, "Supervised Neural Q-learning based Motion Control for Bionic Underwater Robots," Journal of Bionic Engineering, vol. 7, pp. 177-184, Sept. 2010.
- [14] N. Navarro-Guerrero, C. Weber, P. Schroeter, S. Wermter, "Real-world reinforcement learning for autonomous humanoid robot docking", Robotics and Autonomous Systems, vol. 60, pp. 1400-1407. Nov. 2012.
- [15] Su. Shun-Feng, H. Sheng-Hsiung, "Embedding Fuzzy Mechanisms and Knowledge in Box-Type Reinforcement Learning Controllers," IEEE Transaction System, Man, Cybernetic. vol. 32, pp. 645-653, Oct. 2002.