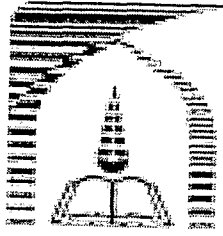


الله الرحمن الرحيم



دانشگاه تربیت مدرس
دانشکده فنی و مهندسی

پایان نامه دوره دکتری مهندسی برق - کنترل

طراحی کنترلر عاملگرای هوشمند برای ناوبری ربات

ولی درهمی

۱۳۸۷ / ۲ / ۵

استاد راهنما:

دکتر وحید جوهری مجد

اساتید مشاور:

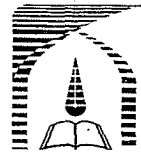
دکتر مجید نیلی احمد آبادی

دکتر حمید رضا مومنی

پاییز ۱۳۸۶

۹۳۲۷۶

کتابخانه تخصصی مهندسی برق
دانشگاه تربیت مدرس



بسمه تعالی

تاییدیه اعضای هیات داوران حاضر در جلسه دفاع از رساله دکتری

آقای ولی درهمی رساله ۳۴ واحدی خود را با عنوان طراحی کنترلگر عامل گرای

هوشمند برای ناوبری ربات در تاریخ ۱۳۸۶/۸/۶ ارائه کردند.

اعضای هیات داوران نسخه نهایی این رساله را از نظر فرم و محتوا تایید کرده و پذیرش

آنها برای تکمیل درجه دکتری مهندسی برق - کنترل پیشنهاد می کنند.

عضو هیات داوران	نام و نام خانوادگی	رتبه علمی	امضا
استاد راهنما	دکتر وحید جوهری مجد	دانشیار	
استاد مشاور	دکتر مجید نیلی احمدآبادی	دانشیار	
استاد مشاور	دکتر حمیدرضا مومنی	دانشیار	
استاد ناظر	دکتر سعید جلیلی	استادیار	
استاد ناظر	دکتر محمد رضا امین ناصری	دانشیار	
استاد ناظر	دکتر نجار اعرابی	دانشیار	
استاد ناظر	دکتر محمد فرخی	دانشیار	
نماینده شورای تحصیلات تکمیلی	دکتر سعید جلیلی	استادیار	

تاییدیه اعضای هیات داوران



دستور العمل حق مالکیت مادی و معنوی در مورد نتایج پژوهشهای علمی دانشگاه تربیت مدرس

مقدمه:

با عنایت به سیاست های پژوهشی دانشگاه در راستای تحقق عدالت و کرامت انسانها که لازمه شکوفایی علمی و فنی است و رعایت حقوق مادی و معنوی دانشگاه و پژوهشگران لازم است اعضای هیأت علمی دانشجویان دانش آموختگان و دیگر همکاران طرح درمورد نتایج پژوهشهای علمی که تحت عناوین پایان نامه رساله و طرحهای تحقیقاتی با هماهنگی دانشگاه انجام شده است موارد ذیل را رعایت نمایند:

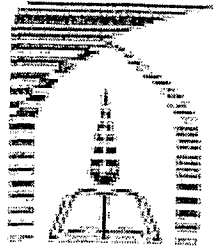
ماده ۱: حقوق مادی و معنوی پایان نامه ها / رساله های مصوب دانشگاه متعلق به دانشگاه است و هر گونه بهره برداری از آن باید با ذکر نام دانشگاه و رعایت آیین نامه ها و دستورالعمل های مصوب دانشگاه باشد.

ماده ۲- انتشار مقاله یا مقالات مستخرج از پایان نامه / رساله به صورت چاپ در نشریات علمی و یا ارائه در مجامع علمی می باید به نام دانشگاه بوده و استاد راهنما نویسنده مسئول مقاله باشند. تبصره: در مقالاتی که پس از دانش آموختگی بصورت ترکیبی از اطلاعات جدید و نتایج حاصل از پایان نامه / رساله نیز منتشر می شود نیز باید نام دانشگاه درج شود.

ماده ۳- انتشار کتاب حاصل از نتایج پایان نامه / رساله و تمامی طرحهای تحقیقاتی دانشگاه باید با مجوز کتبی صادره از طریق حوزه پژوهشی دانشگاه و بر اساس آیین نامه های مصوب انجام می شود. ماده ۴- ثبت اختراع و تدوین دانش فنی و یا ارائه در جشنواره های ملی، منطقه ای و بین المللی که حاصل نتایج مستخرج از پایان نامه / رساله و تمامی طرح های تحقیقاتی دانشگاه باید با هماهنگی استاد راهنما یا مجری طرح از طریق حوزه پژوهشی دانشگاه انجام گیرد.

ماده ۵- این دستورالعمل در ۵ ماده و یک تبصره در تاریخ ۱۳۸۴/۴/۲۵ در شورای پژوهشی دانشگاه به تصویب رسیده و از تاریخ تصویب لازم الاجرا است و هر گونه تخلف از مفاد این دستورالعمل، از طریق مراجع قانونی قابل پیگیری خواهد بود.

نام و نام خانوادگی: ولی دره کمی
تاریخ و امضاء:
۸۶/۱/۱۶



آیین نامه چاپ پایان نامه‌های دانشجویان دانشگاه تربیت مدرس

نظر به اینکه چاپ و انتشار پایان نامه‌های تحصیلی دانشجویان دانشگاه تربیت مدرس، مبین بخشی از فعالیت‌های علمی- پژوهشی دانشگاه است بنابر این بمنظور آگاهی و رعایت حقوق دانشگاه، دانش‌آموختگان این دانشگاه نسبت به رعایت موارد ذیل متعهد می‌شوند:

- ماده ۱: در صورت اقدام به چاپ پایان نامه‌های خود، مراتب را قبلاً بطور کتبی به «دفتر نشر آثار علمی» دانشگاه اطلاع دهد.
- ماده ۲: در صفحه سوم کتاب (پس از برگ شناسنامه)، عبارت ذیل را چاپ کند:
« کتاب حاضر، حاصل رساله دکتری نگارنده در رشته مهندسی برق - کنترل است که در سال ۱۳۸۶ در دانشکده فنی و مهندسی دانشگاه تربیت مدرس به راهنمایی جناب آقای دکتر وحید جوهر مجد و مشاوره جناب آقایان دکتر مجید نیلی احمد آبادی و دکتر حمید رضا مومنی از آن دفاع شده است. »
- ماده ۳: به منظور جبران بخشی از هزینه های انتشارات دانشگاه، تعداد یک درصد شمارگان کتاب (در هر نوبت چاپ) را به « دفتر نشر آثار علمی » دانشگاه اهدا کند. دانشگاه می تواند مازاد نیاز خود را به نفع مرکز نشر در معرض فروش قرار دهد.
- ماده ۴: در صورت عدم رعایت ماده ۳، ۵۰٪ بهای شمارگان چاپ شده را به عنوان خسارت به دانشگاه تربیت مدرس، تادیه کند.
- ماده ۵: دانشجو تعهد و قبول می کند در صورت خودداری از پرداخت بهای خسارت، دانشگاه می تواند خسارت مذکور را از طریق مراجع قضایی مطالبه و وصول کند؛ به علاوه به دانشگاه حق می دهد به منظور استیفای حقوق خود، از طریق دادگاه، معادل وجه مذکور در ماده ۴ را از محل توقیف کتابهای عرضه شده نگارنده برای فروش، تامین نماید.
- ماده ۶: اینجانب ولی درهمی دانشجوی رشته مهندسی برق-کنترل مقطع دکتری تعهدات فوق و ضمانت اجرایی آن را قبول کرده، به آن ملتزم می شوم.

نام و نام خانوادگی: ولی درهمی
تاریخ و امضاء:
۸۶/۱/۱۵

این رساله را تقدیم می‌کنم به

همسر مهربان و فداکارم

و دو فرزند عزیزم

سارا و صدرا

تقدیر و تشکر

خداوند را سپاس می‌گوییم بخاطر لطف و مهربانی هایش و اینکه به من لذت توفیق کسب علم را عطا کرده است و امیدوارم بتوانم تکلیف خود را در برابر آنچه به من ارزانی داشته است بخوبی بانجام برسانم.

از استاد راهنمای محترم جناب آقای دکتر وحید مجد، بخاطر راهنمایی‌های مفیدی که نموده‌اند، و همچنین صبر و حوصله‌ای که در اصلاح مقالات و رساله کشیده‌اند تشکر و سپاسگزاری می‌نمایم.

از استاد مشاور محترم جناب آقای دکتر مجید نیلی احمدآبادی، بخاطر دلسوزی‌ها و راهنمایی‌های ارزنده‌ای که نموده‌اند و وقت زیادی که بسیار بیشتر از آنچه از استاد مشاور انتظار می‌رود برای مشاوره اختصاص دادند، تشکر و قدردانی می‌نمایم. همچنین از جناب آقای دکتر حمید رضا مومنی، استاد مشاور دوم این رساله سپاسگزار هستم.

از زحمات دوست عزیزم آقای دکتر امیر رضا عطارها که در تهیه مقالات انگلیسی کمک بسیار نمودند، قدردانی فراوان می‌نمایم.

در پایان ممنون از صبر، فداکاری، مهربانی، و دلسوزی همسر هستم که اگر حمایت‌های او نبود هرگز موفق به پایان رساندن این مسیر که همراه با فراز و نشیب‌های بسیاری برای من و خانواده‌ام بود، نمی‌شدم.

چکیده

عدم وجود تحلیل ریاضی و عدم کاوش مناسب دو نقطه ضعف عمده در اکثر الگوریتم‌های آموزش تقویتی پیوسته هستند. تحقیقات این رساله بر دو موضوع فوق متمرکز است. یک الگوریتم جدید "آموزش تقویتی فازی" بر مبنای معماری نقاد-تنها، که "آموزش سارسای فازی" (FSL) نامیده می‌شود، ارائه می‌گردد. الگوریتم FSL پارامترهای تالی قواعد سیستم فازی را بصورت روی خط تنظیم می‌کند. انتخاب عمل در هر قاعده بر طبق فرمول جدید ارائه شده در این رساله بنام "بیشینه نرم بهبود یافته" انجام می‌شود. این نحوه انتخاب عمل موجب می‌گردد تا توزیع احتمال انتخاب عمل نهایی در الگوریتم FSL از توزیع بولترمن که یک توزیع پیوسته است، تبعیت کند. وجود نقاط ایستای منطبق بر نقاط ثابت الگوریتم "تکرار تقریب ارزش عمل" (ارائه شده در این رساله) برای FSL اثبات می‌شود. همچنین همگرایی بردار وزن FSL به مقدار یکتا وقتی که از سیاست انتخاب عمل ایستا استفاده شود اثبات می‌گردد. نتایج شبیه سازی در بکارگیری الگوریتم FSL در مقابل الگوریتم "آموزش Q فازی" (FQL)، حاکی از کیفیت و سرعت بالاتر آموزش FSL نسبت به FQL است. همچنین همگرایی FSL در برابر واگرایی FQL در یک مسأله چند هدفه جایکه هر دو الگوریتم از سیاست ایستا استفاده می‌نمایند، نشان داده می‌شود. سپس موضوع تعادل میان "کاوش برای کسب تجربه جدید" و "بهره‌برداری از تجربیات گذشته" در آموزش تقویتی فازی پیوسته مورد بررسی قرار می‌گیرد. کیفیت برقراری تعادل وابستگی شدیدی به دقت تابع ارزش عمل تقریب زده شده دارد. بنابراین، در ابتدا تفاوت‌های بین تخمین تابع ارزش در الگوریتم‌های آموزش تقویتی گسسته با تقریب آن در آموزش تقویتی پیوسته‌ای که از تقریب زنده استفاده می‌کند شرح داده شده و مشکل "فوق برازش" معرفی می‌گردد. برای رفع این مشکل یک روش جدید محاسبه نرخ آموزش ارائه می‌شود. مقدار نرخ آموزش پیشنهادی رابطه معکوس با مقدار ملاقات فازی حالت جاری دارد و باعث تاثیر تقریباً یکسان داده‌های آموزشی روی پارامترهای وزن تقریب زنده و جلوگیری از فوق برازش می‌گردد. برای ایجاد تعادل بین کاوش و بهره‌برداری، یک سیستم فازی T-S مرتبه صفر که مقدار ضریب دما را برای فرمول بیشینه نرم تولید می‌کند، ارائه می‌گردد. با اضافه نمودن دو تکنیک "نرخ آموزش تطبیقی" و "تعادل دهنده فازی کاوش و بهره‌برداری" به الگوریتم FSL، الگوریتم حاصل را FSL بهبود یافته یا EFSL می‌نامیم. نتایج شبیه‌سازی حاکی از برتری قابل توجه در کیفیت و زمان آموزش EFSL در مقابل FSL است. در ادامه راهکارهایی برای چالش‌های موجود در بکارگیری آموزش تقویتی در مسأله ناوبری ربات ارائه می‌شود و ساختار یک کنترلر فازی برای آن تعیین می‌گردد. در انتها عملکرد الگوریتم‌های ارائه شده در تنظیم پارامترهای کنترلر در مسأله مذکور با یکدیگر مقایسه می‌شود.

کلید واژه‌ها: آموزش تقویتی، کاوش، بهره‌برداری، سیستم‌های فازی، ناوبری ربات

فهرست مطالب

صفحه	عنوان
د	فهرست شکل‌ها
و	فهرست جدول‌ها
۱	فصل ۱- مقدمه
۱	۱-۱- پیشگفتار
۲	۲-۱- دو موضوع مهم در آموزش تقویتی پیوسته
۲	۱-۲-۱- تحلیل ریاضی الگوریتم
۲	۲-۲-۱- تعادل بین کاوش و بهره‌برداری از تجربیات
۳	۳-۱- مروری بر مطالعات انجام شده
۳	۱-۳-۱- معماری عملگر- نقاد
۴	۲-۳-۱- معماری نقاد-تنها
۷	۴-۱- هدف و نتایج حاصل از این رساله
۹	۵-۱- ساختار رساله
۱۰	فصل ۲- آموزش تقویتی
۱۰	۱-۲- مقدمه
۱۰	۲-۲- مفاهیم اولیه در RL
۱۲	۳-۲- آموزش تقویتی گسسته
۱۲	۱-۳-۲- آموزش تفاضل موقتی
۱۳	۲-۳-۲- روشهای انتخاب عمل برای RL گسسته
۱۶	۳-۳-۲- تعادل در آموزش تقویتی گسسته
۱۸	۴-۲- آموزش تقویتی پیوسته
۱۹	۱-۴-۲- معماری عملگر- نقاد
۲۰	۲-۴-۲- معماری نقاد-تنها
۲۳	۳-۴-۲- روشهای تعیین نرخ آموزش
۲۳	۴-۴-۲- جمع بندی روشهای آموزش تقویتی پیوسته
۲۴	فصل ۳- آموزش سارسای فازی (FSL)
۲۴	۱-۳- مقدمه
۲۴	۲-۳- آموزش سارسای فازی (FSL)
۲۷	۳-۳- روش تکرار تقریب ارزش عمل (AAVI)
۳۰	۴-۳- تحلیل ریاضی الگوریتم FSL

۳۰	۳-۴-۱- اثبات وجود نقاط ایستا برای الگوریتم FSL
۳۶	۳-۴-۲- همگرایی الگوریتم FSL تحت سیاست ایستا
۳۷	۳-۵- شبیه سازی
۳۷	۳-۵-۱- مسأله هدایت قایق
۴۱	۳-۵-۲- محیط چند هدفه ی بدون مانع
۴۵	۳-۶- جمع بندی و نتیجه گیری
۴۷	فصل ۴- تعادل بین کاوش و بهره‌برداری از تجربیات در آموزش تقویتی
۴۷	۴-۱- مقدمه
۴۷	۴-۲- تعادل در آموزش تقویتی گسسته
۴۷	۴-۲-۱- معرفی چند ایده جدید برای بهبود انتخاب عمل در روش بیشینه نرم
۴۸	۴-۲-۲- شبیه سازی
۵۳	۴-۲-۳- جمع بندی تعادل در آموزش تقویتی گسسته
۵۳	۴-۳- تعادل در الگوریتم‌های FRL با معماری نقاد-تنها
۵۴	۴-۳-۱- وابستگی انتخاب عمل به دقت تقریب تابع ارزش عمل
۵۵	۴-۳-۲- آموزش سارسای فازی بهبود یافته (EFSL)
۶۰	۴-۳-۳- شبیه سازی
۶۷	۴-۳-۴- جزئیات آموزش و نتایج
۶۸	۴-۳-۵- جمع بندی و نتیجه گیری
۷۰	فصل ۵- مسأله ناوبری ربات
۷۰	۵-۱- مقدمه
۷۱	۵-۲- ربات خپرا
۷۳	۵-۳- محیط شبیه ساز KIKS
۷۶	۵-۴- طراحی ساختار کنترلرگر فازی برای مسأله ناوبری
۷۹	۵-۵- شبیه سازی
۸۶	۵-۶- جمع بندی و نتیجه گیری
۸۷	فصل ۶- نتیجه‌گیری و کارهای آینده
۸۷	۶-۱- جمع بندی و نتیجه گیری
۸۸	۶-۲- کارهای آینده
۸۸	۶-۲-۱- اضافه کردن شرایط جدید برای الگوریتم FSL جهت یافتن یک شرط قوی‌تر همگرایی
۸۹	۶-۲-۲- ترکیب آموزش باناظر با الگوریتم FSLE
	۶-۲-۳- استفاده از الگوریتم EFSL در مسائل چند عامله و استفاده از معیارهای خبرگی برای انتشار
۸۹	سیگنال تقویتی
۸۹	۶-۲-۴- استفاده از خطای تفاضل موقتی ارزش-عمل برای تنظیم توابع عضویت مقدم قواعد فازی

- ۵-۲-۶- بررسی وابستگی راندمان به تابع سیگنال تقویتی..... ۸۹
- ۶-۲-۶- مقایسه و تحلیل پدیده سرخوردگی در روانشناسی با مسأله فوق برآزش در آموزش تقویتی پیوسته
..... ۹۰
- ۷-۲-۶- بررسی نحوه مقدار دهی اولیه پارامترهای وزن..... ۹۰
- ۸-۲-۶- بررسی و بکارگیری تکنیک نرخ آموزش تطبیقی در دیگر الگوریتم های آموزشی..... ۹۰
- ۹-۲-۶- تست ایده های ارائه شده در بخش ناوبری ربات برای یک ربات واقعی..... ۹۰
- ۱۰-۲-۶- استفاده از معماری رده بندی در مسأله ناوبری ربات..... ۹۰
- ۹۱..... فهرست مراجع.....
- ۹۵..... واژه نامه فارسی به انگلیسی.....
- ۹۷..... واژه نامه انگلیسی به فارسی.....

فهرست شکل‌ها

صفحه	عنوان
۳	شکل ۱-۱: نمونه ای از به دام افتادن ربات در بین موانع.
۱۰	شکل ۱-۲: چهار چوب آموزش تقویتی.
۱۹	شکل ۲-۲: ساختار عملگر- نقاد.
۳۸	شکل ۱-۳: مسأله هدایت قایق.
۳۸	شکل ۲-۳: توابع عضویت ورودی.
۴۰	شکل ۳-۳: نمونه ای از مسیرهای حرکت قایق.
۴۱	شکل ۴-۳: هیستوگرام <i>LDI</i> ها برای <i>FSL</i> .
۴۱	شکل ۵-۳: هیستوگرام <i>LDI</i> ها در <i>FQL</i> .
۴۳	شکل ۶-۳: محیط چند هدفی بدون مانع.
۴۳	شکل ۷-۳: توابع عضویت ورودی.
۴۴	شکل ۸-۳: مقدارهای وزن در قواعد ۱ و ۲ در <i>FSL</i> .
۴۴	شکل ۹-۳: مقدارهای وزن در قواعد ۳ و ۴ در <i>FSL</i> .
۴۴	شکل ۱۰-۳: منحنی وزن عمل راست در قاعده ۱ برای الگوریتم <i>FQL</i> .
۴۹	شکل ۱-۴: محیط مارپیچ.
۴۹	شکل ۲-۴: مقایسه عملکرد سه روش انتخاب عمل معروف.
۵۰	شکل ۳-۴: مقایسه نتایج عملکرد چهار سیاست پیشنهادی.
۵۱	شکل ۴-۴: مقایسه نتایج عملکرد سیاست بیشینه نرم با سیاست <i>soft-hybrid</i> .
۵۲	شکل ۵-۴: محیط <i>MDP</i> تصادفی.
۵۲	شکل ۶-۴: مقایسه نرخ انتخاب راه بهینه در سیاست بیشینه نرم و چهار سیاست پیشنهادی.
۵۳	شکل ۷-۴: مقایسه نرخ انتخاب راه بهینه در سیاست بیشینه نرم و سیاست ترکیبی.
۵۵	شکل ۸-۴: بلوک دیاگرام الگوریتم <i>EFSL</i> .
۵۹	شکل ۹-۴: توابع عضویت مقدم قواعد تعادل دهنده فازی.
۵۹	شکل ۱۰-۴: خروجی تعادل دهنده فازی برای $E=4.5$.
۶۳	شکل ۱۱-۴: مقدار ملاقات فازی در فضای آموزش برای <i>EFSL</i> .
۶۳	شکل ۱۲-۴: هیستوگرام <i>LDI</i> ها در <i>FSL</i> .
۶۳	شکل ۱۳-۴: هیستوگرام <i>LDI</i> ها در <i>FSL-V</i> .
۶۴	شکل ۱۴-۴: هیستوگرام <i>LDI</i> ها در <i>FSL-A</i> .
۶۴	شکل ۱۵-۴: هیستوگرام <i>LDI</i> ها در <i>EFSL</i> .
۶۴	شکل ۱۶-۴: تعداد دفعات ملاقات هر مربع از فضای آموزش در یک مورد از فوق برازش در <i>FSL</i> .
۶۴	شکل ۱۷-۴: تعداد دفعات ملاقات هر مربع از فضای آموزش در یک اجرای نمونه در <i>EFSL</i> .

- شکل ۴-۱۸: مجموع نرخ های آموزش هر مربع در فضای آموزش در یک مورد از فوق برازش در FSL. ۶۵
- شکل ۴-۱۹: مجموع نرخ های آموزش هر مربع در فضای آموزش در یک اجرای نمونه در EFSL. ۶۵
- شکل ۴-۲۰: ارزش عمل با بالاترین ارزش در FSL. ۶۶
- شکل ۴-۲۱: ارزش عمل با بالاترین ارزش در EFSL. ۶۶
- شکل ۴-۲۲: مسأله آوردن گاری به مرکز. ۶۶
- شکل ۴-۲۳: هیستوگرام LDI ها در FSL. ۶۸
- شکل ۴-۲۴: هیستوگرام LDI ها در EFSL. ۶۸
- شکل ۵-۱: ربات خپرا در مسأله حرکت در کنار دیوار. ۷۲
- شکل ۵-۲: ربات مینیاتوری خپرا. ۷۳
- شکل ۵-۳: ارتباط محیط شبیه سازی با MATLAB. ۷۴
- شکل ۵-۴: مختصات وضعیت ربات و هدف. ۷۵
- شکل ۵-۵: توابع عضویت ورودی های کنترلر. ۷۸
- شکل ۵-۶: محیط آموزشی و نمونه ای از مسیر حرکت ربات در حین آموزش. ۸۰
- شکل ۵-۷: محیط تست برای رویدادهای ۱ تا ۸ در بخش تست. ۸۱
- شکل ۵-۸: محیط تست برای رویداد نهم در بخش تست. ۸۱
- شکل ۵-۹: محیط تست برای رویداد دهم در بخش تست. ۸۲
- شکل ۵-۱۰: هیستوگرام LDI ها در EFSL با $\alpha_0 = 0/1$. ۸۴
- شکل ۵-۱۱: هیستوگرام LDI ها در EFSL با $\alpha_0 = 0/1$. ۸۴
- شکل ۵-۱۲: هیستوگرام LDI ها در FQL با $\alpha_0 = 0/1$. ۸۴
- شکل ۵-۱۳: هیستوگرام LDI ها در FSL با $\alpha_0 = 0/1$. ۸۴
- شکل ۵-۱۴: منحنی تغییرات سیزده پارامتر وزن مربوط به عملهای کاندید در تالی قاعده نوزدهم. ۸۵

فهرست جدول‌ها

صفحه	عنوان
۴۰.....	جدول ۱-۳: نتایج شبیه سازی برای مقادیر متفاوت نرخ آموزش اولیه و ضریب دمای اولیه.....
۵۰.....	جدول ۱-۴: نمایش تعداد انتخاب راه بهینه و متوسط زمان رسیدن به هدف برای سیاست های مختلف.....
۵۲.....	جدول ۲-۴: نرخ رسیدن به هدف بهینه برای سیاست های مختلف.....
۶۲.....	جدول ۳-۴: نتایج شبیه سازی در مسأله هدایت قایق.....
۶۳.....	جدول ۴-۴: مجموع المانهای بردار شدت آتش N برای EFSL.....
۶۸.....	جدول ۵-۴: نتایج شبیه سازی در مسأله آوردن گاری به مرکز.....
۸۳.....	جدول ۱-۵: نتایج شبیه سازی در مسأله ناوبری ربات برای سه الگوریتم.....
۸۵.....	جدول ۲-۵: جزییات نتایج شبیه سازی در مسأله ناوبری ربات برای الگوریتم EFSL.....

فصل ۱- مقدمه

۱-۱- پیشگفتار

طراحی روی خط^۱ کنترلگر بهینه برای سیستم‌های پیچیده در محیط‌هایی که دارای عدم قطعیت^۲ و نایقینی^۳ هستند، اگر غیر ممکن نباشد، کار پیچیده‌ای است. بخصوص برای سیستم‌هایی (مانند ناوبری ربات^۴ ها) که می‌خواهند بطور هوشمند با محیط تعامل داشته باشند. با توجه به اینکه عموماً اطلاعات زیادی در خصوص خروجی مطلوب کنترلگر در این نوع سیستم‌ها وجود ندارد، استفاده از روشهای آموزش بدون ناظر برای طراحی و تنظیم پارامترهای کنترلگر ارجحیت دارد [۱-۳].

آموزش تقویتی (RL)^۵، یک روش قوی مدرن برای آموزش روی خط استراتژی‌های کنترل از طریق تعامل با محیط است. این روش تنها با استفاده از یک معیار اسکالر راندمان، که سیگنال تقویت یا پاداش^۶ نامیده می‌شود، بدون نیاز به سرپرست قادر به آموزش عاملها در محیط‌های پیچیده، ناقطعی و تصادفی می‌باشد [۴-۶]. قابلیت‌های مذکور باضافه قدرت کاوش بالای آموزش تقویتی در جهت یافتن پاسخ بهینه، منجر به استفاده روز افزون آن در وظایف کنترلی پیچیده همچون کنترل حرکت ربات‌ها، و ناوبری ربات شده است.

سازماندهی الگوریتم‌های آموزش تقویتی بر مبنای تخمین ارزش^۷ حالت (یا جفت حالت-عمل) می‌باشد [۶، ۷]. در آموزش تقویتی گسسته مقدار ارزش حالت (یا جفت حالت-عمل) در جدول ارزش ذخیره شده و در هر قدم که آن حالت (یا جفت حالت-عمل) ملاقات شوند به روز رسانی انجام می‌گیرد [۶]. از آنجا که تعداد پارامترهای قابل تنظیم در آموزش تقویتی گسسته، رابطه مستقیمی با عدد اصلی^۸ فضای متغیرهای حالت و عمل مسأله دارد، در مسائل کنترل با فضای حالت-عمل بزرگ و یا پیوسته مانند ناوبری ربات، مشکل تنگنای ابعاد^۹ وجود دارد. بنابراین استفاده از تقریب زنده‌های تابع برای تقریب تابع ارزش^{۱۰} در اینگونه مسائل ضروری است [۸، ۹].

بر این اساس محققین با ترکیب الگوریتم‌های آموزش تقویتی گسسته با تقریب زنده‌هایی همچون شبکه‌های عصبی و منطق فازی، الگوریتم‌های آموزش تقویتی پیوسته را ارائه داده‌اند. کاربردهای گسترده و عملکرد مطلوب سیستم استنتاج فازی در کنترل و مسائل پیچیده و نیز مزایایی چون امکان گنجاندن دانش بشری، ارائه دانش بصورت قواعد اگر-آنگاه، و قابلیت مدلسازی و کنترل سیستم‌های غیر خطی با

¹ Online

² Nondeterministic

³ Uncertainty

⁴ Robot navigation

⁵ Reinforcement Learning

⁶ Reward

⁷ Value

⁸ Cardinality

⁹ Curses of dimensionality

¹⁰ Value function

دقت دلخواه [۱۰-۱۲]، باعث شده است تا محققین [۱۳-۲۶] با ترکیب سیستم‌های فازی بعنوان تقریب زنده با روشهای آموزش تقویتی، الگوریتمهای آموزش تقویتی فازی (FRL) را ارائه دهند. تمرکز ما نیز در این رساله بر روی الگوریتمهای FRL خواهد بود.

توجه شود که در آموزش تقویتی گسسته مقادیر ارزش حالت (یا جفت حالت-عمل) از هم مستقل می‌باشند، در حالیکه در حالت پیوسته به خاطر استفاده از تقریب زنده، مقدار ارزش در هر قدم زمانی بر مبنای پارامترهای تقریب زنده حاصل می‌شود، لذا به روز رسانی پارامترها در هر قدم زمانی بر روی مقادیر ارزش در کل فضا تاثیر می‌گذارد.

۱-۲-۲- دو موضوع مهم در آموزش تقویتی پیوسته

دو موضوع در بکارگیری الگوریتمهای FRL (یا بطور کلی الگوریتمهای آموزش تقویتی پیوسته) در مسائل کنترلی مورد توجه ویژه هستند [۲۷]:

۱-۲-۱- تحلیل ریاضی الگوریتم

بر خلاف نتایج و قضایایی که در خصوص همگرایی مقادیر ارزش برای الگوریتمهای آموزش تقویتی گسسته ارائه شده است [۲۸،۶-۳۰]، اکثر روشهای آموزش تقویتی پیوسته موجود، فاقد تحلیل ریاضی می‌باشند. بنابراین تحلیل ریاضی الگوریتمهای آموزش تقویتی پیوسته، یک موضوع مورد علاقه برای محققین می‌باشد.

۱-۲-۲- تعادل بین کاوش و بهره‌برداری از تجربیات

در حین آموزش تقویتی دو هدف مخالف باید با هم ترکیب شوند: از یک طرف کلیه جفت‌های حالت-عمل مسأله باید به حد کافی برای بدست آوردن دانش مورد کاوش^۲ قرار گیرند و از طرف دیگر تجربه‌های بدست آمده باید در انتخاب عمل بکار گرفته شوند. از نقطه نظر کاوش، مهمترین سؤال این است که چگونه می‌توان زمان یادگیری را حداقل نمود. متناظراً از نقطه نظر بهره‌برداری از تجربیات^۳ گذشته سؤال این است که چگونه می‌توان جریمه‌ها را کاهش داد [۳۱]. برای آموزش مؤثر، عملها باید بگونه‌ای انتخاب شوند که محیط بطور مناسب آزموده شده و از جریمه‌ها نیز حتی الامکان اجتناب گردد. انجام کامل این دو کار بطور همزمان ممکن نیست. در ادامه این رساله برای مختصر نویسی بجای عبارت "تعادل^۴ بین کاوش و بهره‌برداری از تجربیات" از عبارت "تعادل" استفاده می‌کنیم. برقراری تعادل مناسب، ضمن کاهش زمان آموزش، امکان گریز عامل از نقاط اکسترمم محلی را فراهم نموده و منجر به یافتن جواب مطلوب می‌گردد. تعدادی از محققین راهکارهایی برای انتخاب عمل و برقراری تعادل مناسب در آموزش تقویتی گسسته ارائه داده اند، لیکن در خصوص تعادل در آموزش تقویتی پیوسته، نتایج تحقیق جامعی ارائه

¹ Fuzzy Reinforcement Learning

² Exploration

³ Exploitation

⁴ Balance

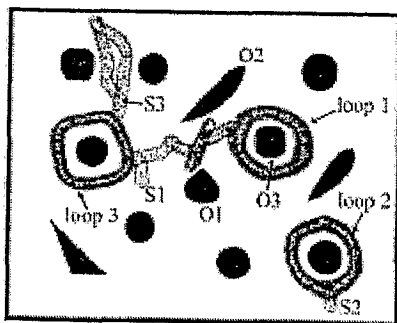
نشده است. بطور کلی مقوله تعادل بین کاوش و بهره‌برداری از تجربیات یک موضوع باز و جدید در آموزش تقویتی پیوسته است که دومین موضوعی است که در این رساله بر روی آن تحقیق می‌شود.

۱-۳-۳- مروری بر مطالعات انجام شده

دو معماری معروف استفاده شده در FRL ها، روش عملگر- نقاد^۱ [۳۲،۱۰] و نقاد- تنها^۲ [۳۳،۶] هستند. در اینجا، ابتدا دو معماری مذکور شرح داده شده و سپس مطالعات انجام شده در خصوص تحلیل ریاضی و تعادل مرور می‌گردد.

۱-۳-۱- معماری عملگر- نقاد

آموزش عملگر- نقاد فازی (FACL)^۳ دارای دو بخش نقاد و عملگر می‌باشد. بخش نقاد برای تقریب تابع ارزش حالت و بخش عملگر برای تولید عمل استفاده می‌گردد. در واقع از بخش نقاد برای تخمین خطای ارزش (خطای تفاضل موقتی (TD)^۴) و بکارگیری این خطا برای تنظیم پارامترهای بخش عملگر استفاده میشود [۱۴]. علیرغم بکارگیری الگوریتمهای FACL در بسیاری از مسائل، هیچ تحلیل ریاضی و یا اثبات همگرایی برای الگوریتمهای FACL موجود که تنها از دو بخش عملگر و نقاد استفاده می‌کنند، موجود نیست. ضعف دیگر ساختار فوق فقدان کاوش مناسب است. در واقع انتخاب عمل در الگوریتم مذکور بر مبنای سیاست حریمانه^۵ می‌باشد و امکان ارزیابی عملهای دیگر در هر حالت را نمی‌دهد. شکل ۱-۱ نمونه‌ای از به دام افتادن ربات را در حلقه‌ی بین موانع حین عبور از آنها در هنگام آموزش با الگوریتم FACL نشان می‌دهد [۲۰].



شکل ۱-۱: نمونه‌ای از به دام افتادن ربات در بین موانع [۲۰].

^۱ Actor-critic
^۲ Critic-only
^۳ Fuzzy Actor Critic Learning
^۴ Temporal Difference
^۵ Greedy

دلیل رخداد مشکل مذکور، فقدان کاوش در الگوریتم‌های FACL می‌باشد. تعدادی از محققین [۳۴،۱۹] برای رفع مشکل فوق، خروجی نهایی را با یک مقدار تصادفی بر مبنای تابع توزیع چگالی احتمال گوسی با میانگین صفر جمع نموده‌اند. علیرغم بهبودهای صورت گرفته به دلیل لحاظ کردن کاوش در کنترلگر، تعادل مناسبی بین کاوش و بهره‌برداری از تجربیات حاصل نشده و نقاط ضعفی چون توزیع احتمال یکنواخت خروجی حول عمل با بالاترین ارزش، و همچنین فقدان ارتباط میان احتمال انتخاب عمل‌ها با مقدارهای ارزش عمل^۱ همچنان پا برجاست.

در [۵] یک الگوریتم FRL جدید با ساختاری شبیه عملگر- نقاد ارائه گردیده و همگرایی روش اثبات شده است. در این الگوریتم بخش نقاد به تقریب مقدار ارزش عمل بجای مقدار ارزش حالت می‌پردازد و خروجی بخش عملگر احتمال انتخاب هر عمل می‌باشد. همچنین بخش دیگری نیز برای تقریب مقدار ارزش در نظر گرفته شده است. گرچه این الگوریتم نسبت به روشهای دیگر FACL دارای تحلیل ریاضی و درجه کاوش بالاتری نسبت به دیگر FACL ها است، لیکن بالاسری محاسباتی بالای آن منجر به عدم گرایش محققین به استفاده از آن شده است. ضمناً به مقوله تعادل بین کاوش و استفاده از تجربیات پرداخته نشده است.

۱-۳-۲- معماری نقاد-تنها

در مقابل معماری عملگر- نقاد، روش نقاد- تنها فقط دارای، یک بخش نقاد است که برای تقریب تابع ارزش عمل استفاده می‌شود و عمل نهایی با توجه به مقادیر ارزش تقریب زده شده تولید می‌گردد. این نحوه انتخاب عمل ضمن فراهم نمودن امکان درجه کاوش بالا، پتانسیل لازم برای برقراری تعادل بین کاوش و بهره‌برداری از تجربیات را دارا می‌باشد. ضمناً ساختار فوق امکان گنجاندن دانش بشری را براحتی ممکن می‌سازد.

در [۲۲-۲۶]، معماری نقاد-تنها با ترکیب یک سیستم فازی با روش آموزش Q^4 (که یک روش "مستقل از سیاست"^۳ است) پیاده سازی شده است. یک نسخه کامل از الگوریتم مذکور با نام آموزش Q فازی (FQL)^۴ در [۲۳] معرفی گردید. در الگوریتم FQL تالی قواعد یک سیستم فازی T-S^۵ مرتبه صفر، بصورت روی خط تنظیم می‌شوند. این الگوریتم در مسائل زیادی از جمله رباتها بکار گرفته شده است [۲۲-۲۶]، اما روش مذکور بر مبنای تجربه بوده و دارای دو ضعف عمده‌ی فقدان تحلیل ریاضی و امکان واگرایی است. مثالها و تحلیل‌هایی از واگرایی الگوریتم FQL و نیز واگرایی الگوریتم‌هایی که از ترکیب آموزش Q با "تقریب زنده‌های تابعی خطی"^۶ (تقریب زنده‌های تابع بصورت ترکیب‌های خطی از

¹ Action value

² Q-Learning

³ Off-policy

⁴ Fuzzy Q-Learning

⁵ Takagi-Sugeno

⁶ Linear function approximators

یک سری توابع پایه [۹]) استفاده می‌کنند در [۳۵،۸] آمده است. در واقع گرچه قابلیت مستقل از سیاست آموزش Q خوشایند است، لیکن تازه سازی مقدار ارزش جفت حالت- عمل بر طبق یک توزیع متفاوت از دینامیک زنجیره مارکف باعث ناپایداری الگوریتم در حالت پیوسته می‌گردد [۳۶]. بعلاوه الگوریتم‌های آموزش Q ترکیب شده با تقریب زنده‌های تابع، ویژگی مستقل از سیاست را از آموزش Q گسسته، به ارث نبرده و مقدار نهایی ارزش جفت حالت- عمل تقریب زده شده به سیاست استفاده شده وابستگی شدید دارد.

بر خلاف نتایج تحلیلی منفی و مثالهایی از واگرایی که برای الگوریتم‌های آموزش تقویتی پیوسته که از تعمیم روش آموزش Q حاصل شده‌اند، موجود است، تعدادی از نتایج تحلیلی مثبت در خصوص ترکیب روشهای آموزش تقویتی "وابسته به سیاست"^۱ با تقریب زنده های تابع وجود دارد [۳۷،۳۹]. وجود نقاط ایستا^۲ در روش تفاضل موقتی خطی که از سیاست انتخاب عمل بیشینه نرم^۳ بهره می‌برد، در [۳۷] اثبات شد. البته این اثبات تنها مربوط به تقریب تابع ارزش حالت است نه تابع ارزش عمل. بعلاوه، با توجه به آنکه در فرمول بیشینه نرم مقدار ارزش عمل مورد نیاز است، در مقاله مذکور مقدار ارزش عمل در هر قدم زمانی، از مقدار تابع ارزش با فرض دانستن مدل محیط محاسبه شده است، در حالیکه در مسائل آموزش تقویتی معمولاً مدل محیط، ناشناخته فرض می‌گردد.

در مقالات [۳۹،۳۸] از ترکیب روش سارسا^۴ (که یک روش وابسته به سیاست است) با تقریب زنده های تابعی خطی (ما این ترکیب را سارسای خطی می‌نامیم) برای تقریب تابع ارزش عمل استفاده شده است. در [۳۸] همگرایی پارامترهای وزن سارسای خطی به یک ناحیه در صورت استفاده از یک سیاست ایستا^۵ در همه رویدادها، اثبات شده است. در [۳۹] همگرایی الگوریتم ارائه شده به پاسخ یکتا به شرط ایستا بودن سیاست ها در هر اپیسود اثبات گردیده است. الگوریتم مذکور دارای دو ضعف عمده می‌باشد: ۱- نتایج قضیه، تضمینی در خصوص کیفیت سیاست نهایی که الگوریتم به آن همگرا می‌شود ارائه نمی‌دهد، ۲- با توجه به آنکه پس از تعیین هر سیاست جدید باید آموزش تا همگرایی وزنها ادامه یابد، سرعت آموزش کند می‌باشد. در حالیکه در مسائل کنترل با آموزش روی خط، مطلوب آن است که در هر قدم زمانی سیاست انتخاب عمل به روز رسانی گردد. توجه شود که هیچگونه پیاده سازی و شبیه سازی برای سه الگوریتم اشاره شده در بالا [۳۷-۳۹] ارائه نشده و در واقع دو چالش عمده برای پیاده سازی این الگوریتم‌ها وجود دارد: انتخاب تابع تقریب زنده تابعی خطی و نحوه انتخاب عمل، به نحوی که این دو بتوانند شرایط بیان شده در لم‌ها و قضایای بیان شده در مقالات مذکور را ارضاء نمایند.

در مجموع، بر خلاف روشهای آموزش تقویتی گسسته، هیچ قضیه‌ای در خصوص همگرایی آموزش تقویتی (فازی) بر مبنای معماری نقاد-تنها وقتی که سیاست انتخاب عمل ایستا نباشد و یا هنگامی که

¹ On-policy

² Stationary points

³ Softmax

⁴ Sarsa

⁵ Stationary

⁶ Episodes

سیاست بر مبنای مقدار ارزش عمل در هر قدم زمانی به روز رسانی شود، وجود ندارد. لذا این موضوع یکی از مواردی است که در این رساله بر روی آن تحقیق می‌گردد.

در ادامه به روشهای انتخاب عمل بکار رفته در الگوریتم‌های اشاره شده در بالا پرداخته می‌شود. الگوریتم‌های FQL ارائه شده در [۲۲-۲۵] از ایده‌های موجود انتخاب عمل و تعادل در حالت گسسته [۶]، [۳۱] برای انتخاب یک عمل در هر قاعده فازی استفاده کرده‌اند. همچنین الگوریتم FQL ارائه شده در [۲۳]، برای انتخاب عملها در هر قاعده از ترکیب فرمول بیشینه نرم با روش "نرخ انجام عمل" [۳۱] بهره برده است. در واقع روشهای انتخاب عمل در مقالات مذکور [۲۲-۲۶] یک گسترش روشهای انتخاب عمل در حالت گسسته بدون توجه به ویژگی‌های تابع ارزش عمل در حالت پیوسته می‌باشد.

در آموزش تقویتی گسسته، تعداد حالت‌ها و عمل‌ها قابل شمارش و محدود می‌باشند. مقادیر حالت‌ها و (یا جفت‌های حالت-عمل) در یک جدول ارزش ذخیره شده و این مقادیر بطور مستقل به روز رسانی می‌گردند. در حالیکه در آموزش تقویتی گسسته، تعداد حالت‌ها و عمل‌ها نامحدود و غیر قابل شمارش می‌باشند. با توجه به استفاده از تقریب زنده‌ها برای تقریب تابع ارزش، به روز رسانی پارامترهای تقریب زنده در هر قدم زمانی می‌تواند باعث تغییر مقدار ارزش در کل فضای مسأله گردد. این تفاوتها باعث می‌شود که ایده‌های ارائه شده برای حالت گسسته چندان در حالت پیوسته کارا نباشند.

نکته دیگر آنکه نحوه انتخاب عمل در مقالات فوق [۲۲-۲۶] هیچ تضمینی برای اینکه تابع توزیع احتمال انتخاب عمل نهایی از یک تابع پیوسته تبعیت کند را ارائه نمی‌دهد. در حالیکه یک نکته مهم در آموزش تقویتی پیوسته، تولید عمل نهایی بر طبق یک توزیع احتمال پیوسته است به گونه‌ای که تغییرات کوچک در مقدار ارزش عمل باعث تغییرات شدید در رفتار عامل نشود [۳۹].

یک مشکل دیگر در RL پیوسته، خطای تقریب تابع ارزش است، بعضاً مقدار این خطا بقدری بزرگ می‌شود که انتخاب عمل بر مبنای آن عملکرد بسیار ضعیفی را باعث می‌گردد. دلیل این خطای زیاد، غالب شدن اثر داده‌های آموزشی ناحیه و یا نواحی می‌باشد که عامل نسبت به نواحی دیگر بسیار بیشتر ملاقات نموده است. این مشکل را ما "فوق برازش"^۱ می‌نامیم و برای اولین بار در این رساله در فصلهای بعد بطور کامل معرفی می‌گردد. برای رفع این معضل ما یک نرخ آموزش تطبیقی را معرفی می‌نماییم. لذا در اینجا اشاره ای هم به نرخ‌های آموزش استفاده شده در FRL ها می‌شود.

بعنوان یک نکته عملی در اکثر مسائل RL، نرخ آموزش با تابعی از زمان کاهش می‌یابد [۶]، [۲۲]، [۲۴]، [۲۵]. از آنجا که نرخ آموزش مذکور تابعی از زمان است، مقدار آن در هنگام به روز رسانی برای پارامترهای قابل تنظیم قواعدی که در شروع آموزش تحریک نمی‌گردند، بقدری کوچک خواهد بود که تغییر این پارامترها محسوس نمی‌باشد. برای رفع مشکل مذکور، در [۴۰] یک نرخ جدید آموزش پیشنهاد گردیده است. در این روش برای پارامترهای قابل تنظیم هر قاعده مستقل از قواعد دیگر یک نرخ آموزش تعیین می‌شود. مقدار نرخ آموزش برای هر قاعده با توجه به مجموع شدت آتش‌های گذشته آن قاعده محاسبه

^۱ Overfitting