





دانشگاه کردستان
دانشکده فنی و مهندسی
گروه مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد رشته مهندسی کامپیوتر گرایش هوش مصنوعی

عنوان:
بهبود کارایی دسته بندی متن بر مبنای ویژگیها و متون
دسته بندی شده مشابه

پژوهشگر:

جلیل بیات

استاد راهنما:

دکتر فریدین اخلاقیان طاب

استاد مشاور:

دکتر پرهام مرادی

مهر ماه ۱۳۹۳

کلیه حقوق مادی و معنوی مترتب بر نتایج مطالعات،

ابتکارات و نوآوری های ناشی از تحقیق موضوع

این پایان نامه (رساله) متعلق به دانشگاه کردستان است.

*** تعهد نامه ***

اینجانب جلیل بیات دانشجوی کارشناسی ارشد رشته مهندسی کامپیوتر گرایش هوش مصنوعی دانشگاه کردستان، دانشکده فنی و مهندسی گروه مهندسی کامپیوتر و فناوری اطلاعات تعهد می‌نمایم که محتوای این پایان نامه نتیجه تلاش و تحقیقات خود بوده و از جایی کپی برداری نشده و به پایان رسانیدن آن نتیجه تلاش و مطالعات مستمر اینجانب و راهنمایی و مشاوره اساتید بوده است.

با تقدیم احترام

جلیل بیات

۱۳۹۳/۶/۲

بسمه تعالی

✽ تعهد نامه دانشجویان تحصیلات تکمیلی دانشگاه کردستان در انجام پایان نامه ✽

(لازم است به عنوان صفحه اول پروپوزال و به عنوان چهارمین برگ پایان نامه و پس از صفحه مشخصات پایان نامه بوده و به دقت

مطالعه و امضا شود)

- | اینجانب | دانشجوی مقطع | رشته | متعهد میشوم: |
|---|--------------|------|--------------|
| ۱- صداقت، امانتداری و بی طرفی را در انجام پژوهش و انتشار نتایج حاصل از آن رعایت نمایم. | | | |
| ۲- در نگارش نتیجه پژوهش های حاصل از موضوع پایان نامه، از باز نویسی نوشته های دیگران بدون ذکر منبع، بازی با الفاظ، زیاده نویسی، کلی گویی و جزم اندیشی و تصرف گرایم پرهیز نمایم و نتایج پژوهشی خود را در موعد مقرر و با اطلاع استاد راهنما منتشر نمایم. | | | |
| ۳- تمامی یافته های مستخرج از پایان نامه متعلق به دانشگاه کردستان بوده و لازم است در کلیه مقالات مستخرج از آنها نام دانشگاه کردستان را تحت عنوان ((دانشجوی دانشگاه کردستان)) یا ((دانش آموخته دانشگاه کردستان)) ذکر نمایم. | | | |
| ۴- در انتشار مقالات نام استاد (استادان) راهنما و استاد (استادان) مشاور را در لیست مولفین مقاله ذکر نمایم و از آوردن اسامی افرادی که نقش موثری در انجام پژوهش نداشته اند، جداً خودداری نمایم. | | | |
| ۵- در بخش سیاستگذاری مقاله، از تمامی افراد و سازمانهایی که در اجرای پژوهش مساعدتی میدول داشته اند با ذکر نوع مشارکت تشکر و قدر دانی نمایم. | | | |
| ۶- از انتشار همپوشان یا ارسال همزمان یک مقاله به چند مجله ویا ارسال مجدد مقاله چاپ شده به مجلات دیگر خودداری نمایم. | | | |
| ۷- در صورت عدم رعایت موارد مذکور، دانشگاه کردستان مجاز خواهد بود تا برابر مقررات اقدام نماید. | | | |

اعضاء و اثر انگشت دانشجو

دستورالعمل نحوه برخورد با موارد تخطی دانشجویان تحصیلات تکمیلی در هنگام انتشار نتایج پژوهش

- در مورد زیر دانشگاه کردستان با مجله مربوطه مکاتبه و درخواست خارج نمودن مقاله را نموده و موضوع را به محل کار یا تحصیل بعدی دانشجو اطلاع خواهد داد.
الف- چاپ مقاله بدون اطلاع و تأیید اسنادان راهنما،
ب- چاپ نتایج حاصل از پژوهش های انجام شده در دانشگاه کردستان بدون ذکر نام دانشگاه
۲- در صورت احراز تخلف از سایر موارد درج شده در تعهد نامه دانشجویی، دانشگاه ضمن مکاتبه با مجله مربوطه، حسب مورد تصمیم گیری خواهد نمود.



دانشگاه کردستان
دانشکده فنی و مهندسی
گروه مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد رشته مهندسی کامپیوتر گرایش هوش مصنوعی

عنوان:

بهبود کارایی دسته بندی متن بر مبنای ویژگیها و متون
دسته بندی شده مشابه

پژوهشگر:

جلیل بیات

در تاریخ / / ۱۳ توسط کمیته تخصصی و هیات داوران زیر مورد بررسی قرار گرفت و با
نمره و درجه به تصویب رسید.

<u>امضاء</u>	<u>مرتبۀ علمی</u>	<u>نام و نام خانوادگی</u>	<u>هیات داوران</u>
	استادیار	دکتر فردین اخلاقیان طاب	۱- استاد راهنما
	استادیار	دکتر پرهام مرادی	۲- استاد مشاور
	استادیار	دکتر محمد فتحی	۳- استاد داور خارجی
	استادیار	دکتر علیرضا عبدالله پوری	۴- استاد داور داخلی

مهر و امضاء معاون آموزشی و تحصیلات تکمیلی دانشکده

مهر و امضاء مدیر گروه

سپاسگذاری و قدردانی

و بعد از مدتها، پس از پیمودن راههای فراوان که با حضور شیرین **اساتید عزیزم**، با راهنماییها و دغدغه‌های فراوانشان، نگاههای **پدر و مادرم**، با چشمهای پر از برق شوق، و زیبایی حضور **خواهرم** در کنارم، که خستگی‌های این راه را به امید و روشنی راه تبدیل کرده و امیدوارم بتوانم در آینده‌ای نزدیک جوابگوی این همه محبت آنها باشم...

اکنون، با احترام فراوان برای این همه تلاش این عزیزان برای موفقیت من....

این پایان نامه را به **پدر و مادرم**، **اساتید عزیز**، **بویژه آقای دکتر فردین اخلاقیان** و **خواهر مهربانم** تقدیم می‌کنم.

امیدوارم قادر به درک زیباییهای وجودشان باشم.

با تشکر

جلیل بیات

چکیده

در یادگیری ماشین، داده‌های آموزشی نقش مهمی را در تعیین کارایی الگوریتم یادگیری ایفا می‌کنند. تعداد این داده‌های آموزشی در طول زمان افزایش می‌یابد و داده‌های جدیدی از راه می‌رسد. این داده‌های جدید، ممکن است دید جدیدی از داده‌های قبلی را ارائه دهند یا توزیع آماری داده‌ها را تغییر دهند. از این رو، فهمیدن اهمیت داده‌های جدید و اجازه به این داده‌ها، برای نقش داشتن در آموزش، به منظور افزایش کارایی سیستم یادگیر، کاری بس زیرکانه است. آموزش دوباره‌ی سیستم یادگیر با همه‌ی داده‌ها، و کنار گذاشتن همه‌ی تجربه‌ی یادگیری گذشته، یکی از راه‌حل‌ها برای این مسئله است ولی این روش دارای نقص‌هایی می‌باشد از جمله اینکه این روش نا کارآمد است و همچنین، این روش قادر به نگه‌داری دانش قبلی نمی‌باشد.

در این پایان‌نامه، ما چهار روش ارائه نمودیم. در روش اول، ما سیستم یادگیری افزایشی برای دسته‌بندی متن ارائه نمودیم، که در این سیستم، برای هر دسته‌ی آموزشی از یک دسته‌بند ماشین بردار پشتیبان برای یادگیری آن استفاده نمودیم و سپس، برای دسته‌بند بعدی علاوه بر دسته‌ی آموزشی آن دسته، از بردارهای پشتیبان دسته‌های آموزشی قبلی به اضافه‌ی داده‌های غلط دسته بندی شده‌ی مجموعه‌های تصدیق دسته‌های آموزشی قبلی استفاده نمودیم. در روش دوم برای اینکه بتوانیم میزان معیارهای ارزیابی کارایی را افزایش و خطای دسته‌بندی را کاهش دهیم به جای استفاده از تنها یک دسته‌بند برای هر دسته‌ی آموزشی، از چند دسته‌بند برای هر دسته‌ی آموزشی استفاده نمودیم و همچنین، به جای در نظر گرفتن وزن یکسان به هر دسته‌بند برای تاثیر در ترکیب خروجی‌ها، از روش اول پایان‌نامه‌ی آقای علی دانش استفاده نمودیم که ایشان به ازای ترکیب هر دسته و دسته‌بند از وزن مستقلی استفاده نموده‌اند. در روش سوم، به منظور کامل کردن روش قبلی، برای هر دسته‌ی آموزشی از میان مجموعه‌ی دسته‌بندها، چندین دسته‌بند قابل قبول را انتخاب نمودیم هدف از ارائه این روش کنار گذاشتن دسته‌بندهای ضعیف و جایگزین کردن آنها با دسته‌بندهای قوی بود که با داده‌های آموزشی بیشتری آموزش دیده بودند بود که این امر اشاره به قابلیت خود تطبیقی سیستم پیشنهادی دارد. در روش چهارم با بسط روش سوم در هنگام گرفتن خروجی نهایی، نودهای ایجاد کننده‌ی نویز را حذف و نودهایی، برای افزایش کارایی دسته‌بندی اضافه نمودیم.

عملکرد روش‌های پیشنهادی با پنج روش دیگر مقایسه شد. نتایج آزمایشات، کارایی روش‌های پیشنهادی و بهبود کارایی دسته‌بندی متن را نشان می‌دهد.

کلمات کلیدی: دسته‌بندی متن، انتخاب ویژگی، یادگیری افزایشی، دسته‌بند، مرحله‌ی آموزشی، دسته‌ی آموزشی

فهرست مطالب

صفحه

عنوان

فصل ۱. مقدمه	۱
۱-۱- مقدمه	۱
۲-۱- بیان مساله تحقیق	۱
۳-۱- اهداف و رویکردها	۲
۴-۱- دستاوردهای پایان نامه	۳
۵-۱- ساختار پایان نامه	۴
فصل ۲. پیشینه	۶
۱-۲- مقدمه	۶
۲-۲- نحوه‌ی بازنمایی اسناد بصورت بردارهای عددی	۷
۱-۲-۲- استخراج ویژگی	۹
۳-۲- دسته‌بندها	۱۲
۱-۳-۲- دسته‌بند ماشین بردار پشتیبان	۱۲
۲-۳-۲- دسته‌بند بیز	۱۵
۳-۳-۲- دسته‌بند درخت تصمیم ID3	۱۶
۴-۳-۲- دسته‌بند نزدیکترین همسایه	۲۰
۵-۳-۲- درخت تصمیم C4.5	۲۱
۴-۲- مروری بر روش‌های قبلی یادگیری افزایشی	۲۲
۱-۴-۲- سطح ۱ تطبیق:	۲۷
۲-۴-۲- سطح ۲ تطبیق: ترکیب دسته‌بندها	۲۷
۳-۴-۲- سطح ۳ از تطبیق: وزندهی به دسته‌بندها	۲۹
۴-۴-۲- یادگیری افزایشی براساس هرس ترکیب دسته‌بندها	۳۴
۵-۲- جمع بندی	۳۸
فصل ۳. روش‌های پیشنهادی	۳۹
۱-۳- مقدمه	۳۹

۲-۳-روش اول: استفاده از داده‌های غلط دسته‌بندی شده و بردارهای پشتیبان دسته‌های آموزشی قبلی	۴۰
۳-۳-روش دوم: اضافه کردن تعداد دسته‌بندها برای هر دسته ی آموزشی و ترکیب دسته بندها براساس روش پیشنهادی آقای علی دانش	۴۳
۳-۴-جمع بندی.....	۴۶
فصل ۴. اصلاح و تکامل روش‌های پیشنهادی	۴۸
۴-۱-مقدمه	۴۸
۴-۲-روش سوم: انتخاب چند دسته‌بند کارا از میان مجموعه‌ی دسته‌بندها و اصلاح دسته‌بندهای با کارایی پایین	۴۸
۴-۳-روش چهارم: حذف و اضافه کردن نود	۵۱
۴-۴-شمای کلی الگوریتم پیشنهادی.....	۵۳
۴-۵-جمع بندی.....	۵۷
فصل ۵. نتایج و تحلیل زمایش	۵۹
۵-۱-مقدمه	۵۹
۵-۲-مجموعه داده‌های مورد استفاده	۵۹
۵-۲-۱-مجموعه‌ی داده‌ی 20Newsgroups.....	۵۹
۵-۲-۲-مجموعه داده‌ی WebKB.....	۶۰
۵-۲-۳-مجموعه داده‌ای Reuters-R8	۶۱
۵-۲-۴-معیارهای میکرو	۶۲
۵-۲-۵-معیارهای ماکرو	۶۳
۵-۳-مقایسه‌ی روش‌های پیشنهادی با یکدیگر	۶۴
۵-۳-۱-بررسی تاثیر افزایش تعداد دسته‌بندها در افزایش کارایی دسته‌بندی	۶۴
۵-۳-۲-بررسی تاثیر روش وزن‌دهی آقای علی دانش به دسته‌بندها در افزایش کارایی دسته‌بندی.....	۶۶
۵-۳-۳-تاثیر استفاده از روش اصلاح نود در بالا بردن کارایی دسته‌بندی.....	۶۸
۵-۳-۴-تاثیر حذف و اضافه نمودن نود در بالا بردن کارایی دسته بندی	۷۰
۵-۴-مقایسه‌ی روش‌های پیشنهادی با سایر روش‌ها	۷۲
۵-۴-۱-مقایسه‌ی روش اول با سایر روش‌ها.....	۷۲

- ۷۴.....مقایسه‌ی روش دوم با سایر روش‌ها ۲-۴-۵
- ۷۷.....مقایسه‌ی روش سوم با سایر روش‌ها ۳-۴-۵
- ۸۰.....مقایسه‌ی روش چهارم با سایر روش‌ها ۴-۴-۵
- ۸۳.....جمع‌بندی ۵-۵
- ۸۴.....فصل ۶. جمع‌بندی و پژوهش‌های آتی.....
- ۸۴.....جمع‌بندی ۱-۶
- ۸۶.....پژوهش‌های آتی ۲-۶

فهرست جداول

صفحه

عنوان

جدول ۱-۵: مشخصات مجموعه داده‌ای ۲۰ گروه خبری.....	۶۰
جدول ۲-۵: مشخصات مجموعه داده‌ای WebKb.....	۶۱
جدول ۳-۵: مشخصات مجموعه داده‌ای Reuters-R8.....	۶۲
جدول ۴-۵: مقادیر اولیه نگه دارنده و مقادیر افزایش آن در هر مرحله.....	۶۴
جدول ۵-۵: مقایسه‌ی روش پیشنهادی اول و روش پیشنهادی دوم برای داده‌های Reuters-R8.....	۶۵
جدول ۶-۵: مقایسه‌ی روش پیشنهادی اول و روش پیشنهادی دوم برای داده‌های WebKb.....	۶۵
جدول ۷-۵: مقایسه‌ی روش پیشنهادی اول و روش پیشنهادی دوم برای داده‌های 20Newsgroup.....	۶۵
جدول ۸-۵: مقایسه‌ی روش پیشنهادی دوم با ترکیب دسته بندها به صورت اکثریت آراء ساده و روش وزن‌دهی روش اول پایان‌نامه‌ی آقای علی دانش برای داده‌های Reuters-R8.....	۶۷
جدول ۹-۵: مقایسه‌ی روش پیشنهادی دوم با ترکیب دسته بندها به صورت اکثریت آراء ساده و روش وزن‌دهی روش اول پایان‌نامه‌ی آقای علی دانش برای داده‌های WebKb.....	۶۷
جدول ۱۰-۵: مقایسه‌ی روش پیشنهادی دوم با ترکیب دسته بندها به صورت اکثریت آراء ساده و روش وزن‌دهی روش اول پایان‌نامه‌ی آقای علی دانش برای داده‌های 20Newsgroup.....	۶۷
جدول ۱۱-۵: مقایسه‌ی روش پیشنهادی دوم با روش پیشنهادی سوم برای مجموعه داده‌ای Reuters-R8.....	۶۸
جدول ۱۲-۵: مقایسه‌ی روش پیشنهادی دوم با روش پیشنهادی سوم برای مجموعه داده‌ای WebKb.....	۶۹
جدول ۱۳-۵: مقایسه‌ی روش پیشنهادی دوم با روش پیشنهادی سوم برای مجموعه داده‌ای 20Newsgroup.....	۶۹
جدول ۱۴-۵: مقایسه‌ی روش پیشنهادی سوم-بخش اول و سوم-بخش دوم با روش پیشنهادی دوم برای مجموعه داده‌ای Reuters-R8.....	۷۰

جدول ۵-۱۵: مقایسه‌ی روش پیشنهادی سوم-بخش اول و سوم-بخش دوم با روش پیشنهادی دوم	
برای مجموعه داده‌ای WebKb.....	۷۱
جدول ۵-۱۶: مقایسه‌ی روش پیشنهادی چهارم-بخش اول و چهارم-بخش دوم با روش پیشنهادی	
سوم برای مجموعه داده‌ای 20Newsgroup.....	۷۱
جدول ۵-۱۷: مقایسه‌ی روش اول با سایر روش‌ها برای داده‌های Reuters-R8.....	۷۳
جدول ۵-۱۸: مقایسه‌ی روش اول با سایر روش‌ها برای داده‌های WebKb.....	۷۳
جدول ۵-۱۹: مقایسه‌ی روش اول با سایر روش‌ها برای داده‌های 20Newsgroup.....	۷۳
جدول ۵-۲۰: مقایسه‌ی روش دوم و سایر روش‌ها برای مجموعه داده‌ای Reuters-R8.....	۷۵
جدول ۵-۲۱: مقایسه‌ی روش دوم با سایر روش‌ها برای مجموعه داده‌ای WebKb.....	۷۶
جدول ۵-۲۲: مقایسه‌ی روش دوم با سایر روش‌ها برای مجموعه داده‌ای 20Newsgroup.....	۷۷
جدول ۵-۲۳: نتایج روش سوم و سایر روش‌ها برای مجموعه داده‌ای Reuters-R8.....	۷۸
جدول ۵-۲۴: نتایج معیارهای ارزیابی روش سوم و سایر روش‌ها برای مجموعه داده‌ای WebKb	۷۹
جدول ۵-۲۵: نتایج معیارهای ارزیابی روش سوم و سایر روش‌ها برای مجموعه داده‌ای	
20Newsgroup.....	۸۰
جدول ۵-۲۶: نتایج معیارهای ارزیابی روش چهارم و سایر روش‌ها برای مجموعه داده‌ای.....	۸۱
جدول ۵-۲۷: نتایج معیارهای ارزیابی روش چهارم و سایر روش‌ها برای مجموعه داده‌ای WebKb	
.....	۸۲
جدول ۵-۲۸: نتایج معیارهای ارزیابی روش چهارم و سایر روش‌ها برای مجموعه داده‌ای	
20NewsGroup.....	۸۲

فهرست شکل ها

صفحه	عنوان
۸.....	شکل ۱-۲: نحوه پیش پردازش متن.....
۱۷.....	شکل ۲-۲: مثالی از یک درخت تصمیم [۳۰].....
۱۸.....	شکل ۳-۲: شبه کد درخت تصمیم ID3 [۳۰].....
۲۰.....	شکل ۴-۲: مثالی از روش نزدیکترین همسایه (K-6) [۱۲].....
۲۲.....	شکل ۵-۲: شبه کد درخت تصمیم C4.5 [۳۴].....
۲۳.....	شکل ۶-۲: دسته‌بندی خودکار متن با استفاده از داده‌های بدون برچسب [۳۵].....
۲۴.....	شکل ۷-۲: شبه کد یادگیری با استفاده از داده‌های بدون برچسب [۳۵].....
۲۵.....	شکل ۸-۲: مدل یادگیری افزایشی با استفاده از ترکیب دسته‌بندها [۱].....
۲۵.....	شکل ۹-۲: شبه کد یادگیری افزایشی با استفاده از ترکیب دسته‌بندها [۱].....
۲۶.....	شکل ۱۰-۲: تطبیق چند سطحی [۳۶].....
۲۸.....	شکل ۱۱-۲: ترکیب یادگیرنده‌های برخط [۳۶].....
۲۹.....	شکل ۱۲-۲: شبه کد وزن‌دهی [۳۶].....
۳۰.....	شکل ۱۳-۲: شبه کد ترکیب اکثریت آراء وزن دار [۳۶].....
۳۱.....	شکل ۱۴-۲: شمای کلی الگوریتم ADAIN [۳۷].....
۳۲.....	شکل ۱۵-۲: شبه کد ADAIN [۳۷].....
۳۳.....	شکل ۱۶-۲: تابع نگاشت با استفاده از MLPs [۳۷].....
۳۴.....	شکل ۱۷-۲: شبه کد هرس [۳۷].....
۳۵.....	شکل ۱۸-۲: شبه کد PBagging++ [۳].....
۳۶.....	شکل ۱۹-۲: دسته‌بندی SV-Incremental [۳۹].....
۳۷.....	شکل ۲۰-۲: شبه کد SV-Incremental [۳۹].....
۴۱.....	شکل ۱-۳: استفاده از بردارهای پشتیبان در بهبود دسته‌بندی [۳۹].....
۴۲.....	شکل ۲-۳: شمای کلی روش اول.....

- شکل ۳-۳: شمای کلی روش دوم..... ۴۵
- شکل ۱-۴: شمای کلی روش سوم..... ۵۰
- شکل ۲-۴: الگوریتم حذف نود ایجاد کننده‌ی نویز..... ۵۲
- شکل ۳-۴: فلوچارت کلی الگوریتم یادگیری افزایشی پیشنهادی برای دسته‌بندی متن..... ۵۶

فصل ۱. مقدمه

۱-۱- مقدمه

مغز انسان دارای توانایی بالایی، در یادگیری افزایشی است. بنابراین شبیه سازی توانایی مغز انسان در یادگیری افزایشی یک از مسائل چالش برانگیز در یادگیری ماشین است. در کاربردهای دنیای واقعی، سه زمینه وجود دارد که نیازمند یادگیری افزایشی است: (۱) امکان جمع آوری داده‌های آموزشی بصورت یکجا امکان پذیر نیست به عنوان مثال در حیطه‌ی پزشکی (۲) در برخی از مسائل جهان واقعی به محض اینکه دسته‌ی آموزشی‌ای رسید ما نیازمند یادگیری هستیم (۳) هنگامی که حجم داده‌های آموزشی زیاد باشد به علت محدود بودن حافظه، قادر به بار کردن آن در حافظه نخواهیم بود [۱].

۱-۲- بیان مساله تحقیق

در یادگیری ماشین، داده‌های آموزشی نقش مهمی را در تعیین کارایی الگوریتم یادگیری ایفا می‌کنند. تعداد این داده‌های آموزشی در طول زمان افزایش می‌یابد و داده‌های جدیدی از راه می‌رسد. این داده‌های جدید، ممکن است دید جدیدی از داده‌های قبلی را ارائه دهند یا توزیع آماری داده‌ها را تغییر دهند. از این رو، فهمیدن اهمیت داده‌های جدید و اجازه به این داده‌ها، برای نقش داشتن در آموزش، به منظور افزایش کارایی سیستم یادگیر، کاری بس زیرکانه است. آموزش دوباره‌ی سیستم یادگیر با همه‌ی داده‌ها، و کنار گذاشتن همه‌ی تجربه‌ی یادگیری گذشته، یکی از راه‌حل‌ها برای این مسئله است ولی این روش دارای نقص‌هایی می‌باشد از جمله اینکه این

روش نا کارآمد است و همچنین، این روش قادر به نگه‌داری دانش قبلی نمی‌باشد. از این رو با توجه به این موضوع، ما بر آن شدیم که سیستمی را ارائه نماییم که قادر باشد دانش یاد گرفته شده‌ی قبلی را حفظ نماید در عین حال، قادر به یادگیری از داده‌های جدید نیز باشد. چنین سیستمی، علاوه بر این مزیت‌ها باید بتواند از دانش یاد گرفته شده‌ی قبلی در دسته‌بندی متون فعلی استفاده کند و علاوه بر آن باید بتواند دقت بالایی نیز در دسته‌بندی متون داشته باشد.

۱-۳- اهداف و رویکردها

با توجه به مطالب ذکر شده در مورد اهمیت یادگیری افزایشی، هدف ما در این پایان‌نامه ارائه سیستمی بود که قادر باشد علاوه بر حفظ دانش قبلی داده‌های جدید را نیز فرا گیرد از این رو چنین سیستمی بایستی دارای قابلیت‌های زیر باشد:

- چنین سیستمی بایستی قابلیت بسط و گسترش داشته باشد؛ بدین معنی که قابلیت بالایی در یادآوری داده‌های یاد گرفته شده‌ی قبلی داشته باشد در عین حال بتواند در هنگام یادگیری داده‌های جدید، داده‌هایی از داده‌های یاد گرفته شده‌ی قبلی را به نحوی در این آموزش دخالت دهد که خط تصمیم ساخته شده در هر مرحله، به خط بهینه‌ی سراسری نزدیک شود.
- چنین سیستمی باید قادر باشد نقاط ضعف دانش یاد گرفته‌ی قبلی را بهبود بخشد و این نقاط ضعف را پوشش دهد.
- چنین سیستمی باید بتواند قابلیت خود تطبیقی داشته باشد بدین معنی که در حین آموزش فرضیه‌های ضعیف را شناسایی کرده و آنها را با فرضیه‌های قویتر جایگزین نماید.
- این سیستم باید بتواند بعد از اتمام مراحل آموزشی و در مرحله‌ی تست نهایی فرضیه‌های ایجاد کننده‌ی نویز را شناسایی نماید.
- این سیستم باید در مرحله‌ی تست نهایی با توجه به اینکه همه‌ی فضای داده‌ای را در اختیار دارد باید قادر باشد تا جایی که امکان دارد کارایی سیستم را افزایش دهد.
- دارای کارایی بالایی نسبت به روش‌های ارائه شده‌ی قبلی باشد.

۱-۴- دستاوردهای پایان نامه

در راستای رسیدن به اهداف پایان نامه چهار روش بصورت زیر ارائه شد:

در روش اول، ما سیستم یادگیری افزایشی ارائه نمودیم که در این سیستم، برای هر دسته‌ی آموزشی از یک دسته‌بند ماشین بردار پشتیبان برای یادگیری آن استفاده نمودیم. سپس، برای دسته‌ی آموزشی بعدی علاوه بر دسته‌ی آموزشی آن دسته، از بردارهای پشتیبان دسته‌های آموزشی قبلی به اضافه‌ی داده‌های غلط دسته‌بندی شده‌ی مجموعه‌های تصدیق دسته‌های آموزشی قبلی استفاده نمودیم. سپس، این دسته‌بندهای آموزش دیده در هر مرحله را ذخیره نمودیم. در پایان به منظور ارزیابی کارایی سیستم پیشنهادی یک مجموعه‌ی تست نهایی را به هر کدام از دسته‌بندهای آموزش دیده ارائه نمودیم و برای تولید خروجی نهایی، از روش اکثریت آراء ساده استفاده نمودیم. استفاده از بردارهای پشتیبان دسته‌های آموزشی قبلی، به منظور استفاده از تجربیات گذشته‌ی دسته‌بندهای پیشین در دسته‌بندهای کنونی و ساخت خط تصمیم بهینه بود و همچنین، استفاده از داده‌های غلط دسته‌بندی شده، برطرف کردن نقاط ضعف دسته‌بندهای پیشین و ساخت دسته‌بندهای قوی بود.

در روش دوم برای اینکه بتوانیم میزان معیارهای ارزیابی کارایی را افزایش و خطای دسته‌بندی را کاهش دهیم به جای استفاده از تنها یک دسته‌بند برای هر دسته‌ی آموزشی، از چند دسته‌بند برای هر دسته‌ی آموزشی استفاده نمودیم. همچنین، به جای در نظر گرفتن وزن یکسان به هر دسته‌بند برای تاثیر در ترکیب خروجی‌ها، از روش اول پایان‌نامه‌ی آقای علی دانش استفاده نمودیم که ایشان به ازای ترکیب هر دسته و دسته‌بند از وزن مستقلی استفاده نموده‌اند و به منظور بروزرسانی مداوم این وزن‌ها در هر مرحله از یادگیری افزایشی، مجموعه تصدیق فعلی به اضافه‌ی تمام مجموعه تصدیق‌های قبلی را به همه‌ی دسته‌بندهای فعلی و همه‌ی دسته‌بندهای قبلی ارائه نمودیم که این امر سبب ثبات نسبی این وزن‌ها گردید.

در روش سوم، به منظور کامل کردن روش قبلی سیستم یادگیری افزایشی ارائه نمودیم که قادر بود برای هر دسته‌ی آموزشی از میان مجموعه‌ی دسته‌بندا، چندین دسته‌بند قابل قبول را انتخاب نماید. هدف از ارائه این روش کنار گذاشتن دسته‌بندهای ضعیف و جایگزین کردن آنها با دسته‌بندهای قوی‌تری که با داده‌های آموزشی بیشتری آموزش دیده بودند بود که این امر اشاره به قابلیت خود تطبیقی سیستم دارد.

در روش چهارم با بسط روش سوم در هنگام گرفتن خروجی نهایی، نودهایی را حذف و نودهایی را اضافه نمودیم. هدف از ارائه این روش حذف نودهای ایجاد کننده‌ی نویز که باعث پایین آمدن کارایی سیستم می‌شد بود و اضافه کردن نودهایی بود که دارای کارایی بالایی بودند که این امر سبب افزایش کارایی کلی سیستم یادگیر تا جایی که امکان داشت می‌شد.

۱-۵- ساختار پایان نامه

این پایان نامه در شش فصل تنظیم گردیده است که رئوس مطالب و موضوعات اصلی هر فصل در ذیل بطور خلاصه بیان گردیده است:

• **فصل اول** به انگیزه، اهمیت و شرح موضوع تحقیق می‌پردازد. همچنین ساختار کلی پایان نامه و خلاصه مطالب ارائه شده در این تحقیق نیز بیان می‌گردد.

• **فصل دوم** در این فصل ابتدا به مطالعه‌ی نحوه بازنمایی اسناد بصورت بردارهای عددی می‌پردازیم سپس، به مطالعه‌ی انواع روش‌های انتخاب و ویژگی فیلتر می‌پردازیم در ادامه، به بررسی دسته‌بندی‌های ارائه شده در این پایان‌نامه می‌پردازیم و در آخر به مرور کارهای ارائه شده‌ی قبلی در این مورد می‌پردازیم.

• **فصل سوم** ما در این فصل به ارائه‌ی دو روش می‌پردازیم. در روش اول ما سیستم یادگیری افزایشی ارائه نمودیم که در این سیستم، برای هر دسته‌ی آموزشی از یک دسته‌بند ماشین بردار پشتیبان برای یادگیری آن استفاده نمودیم و سپس، برای دسته‌ی آموزشی بعدی علاوه بر دسته‌ی آموزشی آن دسته، از بردارهای پشتیبان دسته‌های آموزشی قبلی به اضافه‌ی داده‌های غلط دسته بندی شده‌ی مجموعه‌های تصدیق دسته‌های آموزشی قبلی استفاده نمودیم. سپس، این دسته‌بندی‌های آموزش دیده در هر مرحله را ذخیره نمودیم و در پایان به منظور ارزیابی کارایی سیستم پیشنهادی یک مجموعه‌ی تست نهایی را به هر کدام از دسته‌بندی‌های آموزش دیده ارائه نمودیم و برای تولید خروجی نهایی، از روش اکثریت آراء ساده استفاده نمودیم. در روش دوم برای اینکه بتوانیم میزان معیارهای ارزیابی کارایی را افزایش و خطای دسته‌بندی را کاهش دهیم به جای استفاده از تنها یک دسته‌بند برای هر دسته‌ی آموزشی، از چند دسته‌بند برای هر دسته‌ی آموزشی استفاده نمودیم و همچنین، به

جای در نظر گرفتن وزن یکسان به هر دسته‌بند برای تاثیر در ترکیب خروجی‌ها، از روش اول پایان‌نامه‌ی آقای علی دانش استفاده نمودیم که ایشان به ازای ترکیب هر دسته و دسته‌بند از وزن مستقلی استفاده نموده‌اند و به منظور بروزرسانی مداوم این وزن‌ها در هر مرحله از یادگیری افزایشی، مجموعه تصدیق فعلی به اضافه‌ی تمام مجموعه تصدیق‌های قبلی را به همه‌ی دسته‌بندهای فعلی و همه‌ی دسته‌بندهای قبلی ارائه نمودیم

• **فصل چهارم** در این فصل ما سعی در بهبود دو روش ارائه شده در فصل قبل پرداختیم و بدین ترتیب دو روش جدید مطرح نمودیم. در روش سوم، به منظور کامل کردن روش قبلی، سیستم یادگیری افزایشی ارائه نمودیم که قادر بود برای هر دسته‌ی آموزشی از میان مجموعه‌ی دسته‌بندها، چندین دسته‌بند قابل قبول را انتخاب نماید. در روش چهارم با بسط روش سوم در هنگام گرفتن خروجی نهایی، نودهایی را حذف و نودهایی را اضافه نمودیم.

• **فصل پنجم** این فصل از پنج قسمت اصلی تشکیل شده است دربخش اول به معرفی مجموعه‌های داده‌ای می‌پردازیم سپس، در بخش دوم به معرفی معیارهای ارزیابی پرداخته و در بخش سوم و چهارم به ترتیب به مقایسه‌ی روش‌های پیشنهادی با یکدیگر و به مقایسه‌ی روش‌های پیشنهادی با روش‌های قبلی می‌پردازیم و سپس، در بخش آخر به تحلیل وزن دسته‌بندها که براساس روش اول پایان‌نامه‌ی آقای علی دانش وزن‌دهی شده‌اند می‌پردازیم.

• **فصل ششم** در این فصل راهکارهای ارائه شده در این پایان‌نامه بطور خلاصه بررسی شده و در ادامه پیشنهاداتی برای توسعه و بهینه‌سازی هر چه بیشتر این سیستم ارائه می‌شود.