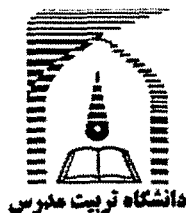


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ
الْحَمْدُ لِلَّهِ الَّذِي
خَلَقَ السَّمَوَاتِ وَالْأَرْضَ
وَالَّذِي يُرِيهِمْ آيَاتِهِ
وَالَّذِي يُخْرِجُ النَّوْمَ
وَالَّذِي يُخْرِجُ النَّوْمَ
وَالَّذِي يُخْرِجُ النَّوْمَ

۹۷۱۹۵



دانشگاه تربیت مدرس

دانشکده فنی و مهندسی

پایان‌نامه‌ی دوره کارشناسی ارشد مهندسی کامپیوتر - نرم‌افزار

ارتقای کیفیت انتخاب خصیصه متون با روش‌های الگوریتم ژنتیک و بهینه‌سازی گروه ذرات

شیما گرانی

استاد راهنما

دکتر سعید جلیلی

وزارت معارف و اوقاف و صنایع مستظرفه

۱۳۸۷ / ۲ / ۵

زمستان ۱۳۸۶

۹۳۱۹۵



بسمه تعالی

تاییدیه اعضای هیات داوران حاضر در جلسه دفاع از پایان

خانم شیما گرانی پایان نامه واحدی خود را با عنوان ارتقای کیفیت انتخاب خصیصه

متون با روشهای الگوریتم ژنتیک و بهینه سازی گروه ذرات در تاریخ

۱۳۸۶/۱۱/۳ ارائه کردند.

اعضای هیات داوران نسخه نهایی این پایان نامه را از نظر فرم و محتوا تایید کرده و پذیرش آنرا برای تکمیل درجه کارشناسی ارشد مهندسی برق - مهندسی کامپیوتر نرم افزار پیشنهاد می کنند.

امضا	رتبه علمی	نام و نام خانوادگی	عضو هیات داوران
	استادیار	دکتر سعید جلیلی	استاد راهنما
	استاد	دکتر احسان اله کبیر	استاد ناظر
	استادیار	دکتر نصراله مقدم چرکری	استاد ناظر
	دانشیار	دکتر جعفر حبیبی	استاد ناظر
	استادیار	دکتر نصراله مقدم چرکری	مدیر گروه (یا نماینده گروه تخصصی)

این نسخه به عنوان نسخه نهایی پایان نامه / رساله مورد تایید است.

امضای استاد راهنما:

آیین نامه چاپ پایان نامه (رساله) های دانشجویان دانشگاه تربیت مدرس

نظر به بنکة چاپ و نسنار بیان نامه (رساله) های تحصینی دانشجویان دانشگاه تربیت مدرس، مبین بخشی از فعالیت های علمی - پژوهشی دانشگاه است بنابراین به منظور آگاهی و رعایت حقوق دانشگاه، دانش آموختگان این دانشگاه نسبت به رعایت موارد ذیل متعهد می شوند:

ماده ۱: در صورت اقدام به چاپ پایان نامه (رساله) ی خود، مراتب را قبلاً به طور 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 پایان نامه کارشناسی ارشد/ رساله دکتری نگارنده در رشته مهندسی کامپیوتر است که در سال ۱۳۸۶ در دانشکده فنی دانشگاه تربیت مدرس به راهنمایی سرکار خانم/جناب آقای دکتر سعید جلیلی، مشاوره سرکار خانم/جناب آقای دکتر

و مشاوره سرکار خانم/جناب آقای دکتر 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 رابه عنوان خسارت به دانشگاه تربیت مدرس، تأدیه 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 سید رانی مقطع کارشناسی ارشد دانشجوی رشته مهندسی کامپیوتر نرم افزار

تعهد فوق و ضمانت اجرایی آن را قبول کرده، به آن ملتزم می شوم.

نام و نام خانوادگی: سید رانی

تاریخ و امضا:

۸۶/۱۱/۲۰

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

۱۴۹۱۴

مقدمه: با عنایت به سیاست‌های 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 هیات علمی، دانشجویان، دانش‌آموختگان و دیگر همکاران طرح، در مورد نتایج پژوهشهای علمی که تحت عناوین پایان‌نامه، رساله و طرحهای تحقیقاتی که با هماهنگی دانشگاه انجام شده است، موارد ذیل را رعایت نمایند:

ماده ۱- حقوق مادی و معنوی پایان‌نامه‌ها / رساله‌های مصوب دانشگاه متعلق به دانشگاه است و هرگونه بهره‌برداری از آن نباید با ذکر نام دانشگاه و رعایت آیین‌نامه‌ها و دستورالعمل‌های مصوب دانشگاه باشد.

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

ماده ۳- انتشار کتاب حاصل از نتایج پایان‌نامه / رساله و تمامی طرحهای تحقیقاتی دانشگاه باید با مجوز کتبی صادره از طریق حوزه پژوهشی دانشگاه و بر اساس آئین‌نامه‌های مصوب انجام می‌شود.

ماده ۴- ثبت اختراع و تدوین دانش فنی و یا ارائه در جشنواره‌های ملی، منطقه‌ای و بین‌المللی که حاصل نتایج مستخرج از پایان‌نامه / رساله و تمامی طرحهای 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 ۱۳۸۲/۴/۲۵ در شورای پژوهشی دانشگاه به تصویب رسیده و از تاریخ تصویب لازم الاجرا است و هرگونه تخلف از مفاد این دستورالعمل، از طریق مراجع قانونی قابل پیگیری خواهد بود.

تقدیم به پدر و مادر عزیزم

برادر مهربانم

و همسر عزیزم

تشکر و قدردانی

سپاس خداوند منان را که بار دیگر دری از درهای رحمت خود را بر من گشود و به من توفیق کسب علم، در راستای افزایش درک خویش از عظمت باری تعالی و شکوه جهان هستی مرحمت نمود. زیرا هر آنچه که بر انسان اتفاق می‌افتد بهانه‌ای است برای شناخت بهتر از جهان هستی و نزدیکی بیشتر با باری تعالی؛ این پژوهش نیز از این قاعده مستثنی نیست هر چند که در نگاه اول موضوع پژوهش چنین حسی را تداعی نمی‌کند.

در فرایند انجام این پژوهش از استاد گرانقدر جناب آقای دکتر سعید جلیلی درس‌های علمی و اخلاقی بی‌شماری را آموختم. از ایشان به دلیل این که طی مدت زمان انجام این پژوهش تجربیات ارزنده و چندین ساله خود را در اختیار اینجانب قرار دادند، نهایت تشکر را دارم. از جناب آقای دکتر نصراله مقدم، به دلیل حمایت‌های بی‌دریغ‌شان و اعضای محترم هیات داوران به دلیل صرف وقت برای ارزیابی این پژوهش، نهایت تشکر را دارم.

آخرین و نه کمترین سپاس من از آقایان محسن ملانوری، مهدی آبادی و محرم منصوری زاده و خانمها پریسا جلیلی، پرواز مهدابی و نرگس خاکپور به دلیل صرف وقت و راهنمایی‌هایشان است.

لازم به ذکر است که این پروژه با استفاده از حمایت مالی مرکز تحقیقات مخابرات ایران انجام شده است.

شیما گرانی

۸۶/۱۱/۲

چکیده

در دسته‌بندی متون، بطور معمول از کلمات متن به‌عنوان خصیصه‌های متن استفاده می‌شود. با توجه به آنکه هر مجموعه متون دارای تعداد زیادی کلمه است، روش‌های دسته‌بندی متون با تعداد زیادی خصیصه مواجه هستند. این تحقیق، دو فرایند و چهار روش جهت بهبود انتخاب خصیصه در حوزه متن پیشنهاد می‌دهد و به ارزیابی و مقایسه روش‌های پیشنهادی با یکدیگر و با سایر روش‌های موجود می‌پردازد.

در فرایند اول (SFSP)، ابتدا با استفاده از کل خصیصه‌ها دسته‌بند یادگیری شده و کارایی آن محاسبه می‌شود. سپس خصیصه‌هایی که باعث کاهش کارایی شده‌اند و در ادامه خصیصه‌هایی که تأثیری در کارایی دسته‌بند ندارد، حذف می‌گردند. در نهایت اگر حذف بیشتر خصیصه‌ها مطلوب باشد، حذف خصیصه‌ها با کارایی کمتر ولی قابل قبول ادامه می‌یابد. برای پیاده‌سازی فرایند SFSP دو روش پیشنهاد شده است که در روش اول (SGA) از الگوریتم ژنتیک و در روش دوم (SBPSO) از BPSO برای جستجوی فضای خصیصه‌ها استفاده شده است. برای ارزیابی زیرمجموعه خصیصه‌های انتخاب شده، از روش دسته‌بندی SVM در هر دو روش استفاده شده است. روش SGA، با کارایی $0/9676$ حدود $98/83$ درصد خصیصه‌ها و روش SBPSO با کارایی $0/9684$ حدود $85/34$ درصد خصیصه‌ها را از متون روترز-21578 گونه ModApte حذف می‌نماید.

در فرایند پیشنهادی دوم، MFSP، برای فرار از بهینه محلی، در شرایطی که ممکن است GA و یا PSO در بهینه محلی گیر کنند، بعد از تعدادی تکرار که جمعیت تا حدی به سمت یک راه حل همگرا شد، خصیصه‌هایی که تا این لحظه به عنوان نویز شناخته شده‌اند دور ریخته می‌شوند و در فضای خصیصه‌های باقی مانده به دنبال حذف نویزهای بیشتر می‌گردیم. به این ترتیب، کروموزوم‌ها در GA و یا ذرات در PSO با توجه به بردار خصیصه‌های جدید، دوباره به صورت تصادفی مقداردهی می‌شوند. شرط پایان برای الگوریتم کاهش خصیصه‌ها در فرایند MFSP، آن است که در مرحله جدید قادر به افزایش کارایی نسبت به مرحله قبل نباشیم. پیاده‌سازی MFSP با استفاده از GA به عنوان روش جستجو در روش پیشنهادی سوم، MGA و با استفاده از BPSO به عنوان روش جستجو در روش پیشنهادی چهارم، MBPSO، انجام شده است. در این دو روش نیز، برای ارزیابی زیرمجموعه خصیصه‌های انتخاب شده، از روش دسته‌بندی SVM استفاده شده است. ارزیابی روش‌های MGA و MBPSO نشان داد که روش MGA، $99/05$ درصد خصیصه‌ها را با کارایی $0/972$ و روش MBPSO نیز $98/5$ درصد خصیصه‌ها را با کارایی $0/98$ حذف نمود. مقایسه روش‌های پیشنهادی با IG و CHI به عنوان دو نمونه از بهترین روشهای فیلتری، برتری روشهای پیشنهادی را نشان می‌دهد.

برای سنجش کارایی روش‌های پیشنهادی در سایر دامنه‌ها، روش SGA در انتخاب خصیصه هرزنامه‌های متنی و همچنین روش‌های SGA و MBPSO در زمینه انتخاب ژن برای دسته‌بندی نسوج تومور به‌کار گرفته شد. مقایسه روش SGA با سایر روش‌های جداسازی هرزنامه‌ها، نشان داد که روش SGA به خوبی قابل مقایسه با بهترین روش‌های آماری موجود است. اعمال روش SGA و MBPSO بر روی مجموعه داده سرطان روده بزرگ برای جداسازی نسوج تومور نشان داد که این روش‌ها می‌توانند با انتخاب تعداد بسیار کم (۳ یا ۴) ژن کارایی دسته‌بندی تومور را افزایش دهند. روش‌های پیشنهادی می‌توانند با شناساندن ژن‌های مهم در دسته‌بندی نسوج تومور، در سایر جنبه‌های پزشکی مفید واقع شوند.

در زمینه دسته‌بندی متون با روش‌های تکاملی، روش Liu و همکاران، بهبود داده شده به طوری که روش دسته‌بندی پیشنهادی، از نظر دقت، فراخوانی، F1 و صحت به ترتیب ۲، ۱۵، ۹ و ۲ درصد بهتر از روش Liu و همکارانش عمل کرده است.

کلمات کلیدی: دسته‌بندی متون، انتخاب خصیصه، روش‌های تکاملی، روش بهینه‌سازی گروه ذرات

فهرست مطالب

عنوان.....	صفحه.....
۱- کلیات.....	۱
۱-۱- مقدمه.....	۱
۲-۱- اهداف و نتایج پژوهش.....	۲
۳-۱- مروری بر ساختار پایان نامه.....	۵
۲- مفاهیم پایه در دسته‌بندی متون.....	۷
۱-۲- مقدمه.....	۷
۲-۲- متن و خصیصه‌های آن.....	۷
۳-۲- مرحله پیش پردازش.....	۸
۴-۲- نمایش متون.....	۸
۵-۲- نمایش معنایی.....	۱۱
۶-۲- مجموعه آموزش، مجموعه آزمایش، مجموعه آزمایش صحت.....	۱۲
۷-۲- معیارهای ارزیابی کارایی.....	۱۳
۸-۲- روش‌های میانگین‌گیری جزئی و میانگین‌گیری کلی.....	۱۵
۹-۲- مجموعه داده‌های استاندارد برای ارزیابی دسته‌بندی متون.....	۱۶
۱۰-۲- نتیجه‌گیری.....	۱۷
۳- تاریخچه پژوهش.....	۱۸
۱-۳- مقدمه.....	۱۸
۲-۳- روش‌های انتخاب خصیصه.....	۱۹
۳-۳- روش‌های آماری انتخاب خصیصه متون.....	۲۰
۴-۳- روش‌های محاسباتی بر گرفته از طبیعت در انتخاب خصیصه.....	۲۳
۵-۳- روش‌های دسته‌بندی متون.....	۲۶
۶-۳- نتیجه‌گیری.....	۳۱
۴- روش‌های محاسباتی بر گرفته از طبیعت.....	۳۲
۱-۴- مقدمه.....	۳۲
۲-۴- الگوریتم ژنتیک.....	۳۲

۳۴	۳-۴- بهینه‌سازی گروه ذرات
۴۲	۴-۴- نتیجه گیری
۴۳	۵- روش‌های پیشنهادی انتخاب خصیصه
۴۳	۱-۵- مقدمه
۴۳	۲-۵- اولین فرایند پیشنهادی انتخاب خصیصه (SFSP)
۴۵	۳-۵- روش اول- پیاده سازی SFSP با استفاده از الگوریتم ژنتیک (SGA)
۵۰	۴-۵- روش دوم- پیاده سازی SFSP با استفاده از BPSO (SBPSO)
۵۳	۵-۵- دومین فرایند پیشنهادی انتخاب خصیصه (MFSP)
۵۳	۶-۵- روش سوم- پیاده سازی MFSP با استفاده از الگوریتم ژنتیک (MGA)
۵۷	۷-۵- روش چهارم- پیاده سازی MFSP با استفاده از BPSO (MBPSO)
۵۸	۶- ارزیابی روش‌های پیشنهادی انتخاب خصیصه در حوزه متون
۵۸	۱-۶- مجموعه داده
۵۹	۲-۶- ارزیابی روش اول (SGA)
۶۴	۳-۶- ارزیابی روش دوم (SBPSO)
۶۶	۴-۶- ارزیابی روش سوم (MGA)
۶۹	۵-۶- ارزیابی روش چهارم (MBPSO)
۷۳	۶-۶- مقایسه روش‌های پیشنهادی با یکدیگر
۷۵	۷-۶- مقایسه روش‌های پیشنهادی با سایر روش‌های انتخاب خصیصه
۷۶	۸-۶- نتیجه گیری
۷۷	۷- کاربرد روش‌های پیشنهادی در سایر دامنه‌ها
۷۷	۱-۷- جداسازی هرز نامه‌های متنی
۸۵	۲-۷- انتخاب ژن در دسته‌بندی نسوج تومور
۹۴	۸- دسته‌بندی متون با رویکرد الگوریتم ژنتیک
۹۴	۱-۸- پژوهش‌های مرتبط
۹۵	۲-۸- انتخاب خصیصه و تولید بردار خصیصه در روش دسته‌بندی پیشنهادی
۹۶	۳-۸- آموزش دسته‌بند با الگوریتم ژنتیک
۹۷	۴-۸- ارزیابی روش پیشنهادی
۱۰۱	۵-۸- نتیجه‌گیری

۱۰۲ نتیجه‌گیری و پیشنهادات	۹-۱
۱۰۲ مقدمه	۹-۱
۱۰۳ نتیجه‌گیری	۹-۲
۱۰۵ پیشنهادات	۹-۳
۱۰۷ فهرست اختصارات	
۱۰۹ فهرست واژگان انگلیسی- فارسی	
۱۱۲ فهرست واژگان فارسی- انگلیسی	
۱۱۵ مراجع و منابع	

فهرست جداول

عنوان.....	صفحه.....
جدول ۱-۲- روش نمایش برداری اسناد متنی.....	۸.....
جدول ۱-۶- آمار سندهای Reuters-21578(10) به تفکیک ۱۰ کلاس.....	۵۸.....
جدول ۲-۶- آمار سندهای Reuters-21578 برای زیرمجموعه انتخابی.....	۵۹.....
جدول ۳-۶- کارایی دسته‌بند روی داده‌های MODAPTE رویترز-۲۱۵۷۸ با تمامی خصیصه‌ها.....	۶۰.....
جدول ۴-۶- تغییرات کارایی روش SGA به تفکیک مراحل و درصد‌های باقیمانده از خصیصه‌ها.....	۶۲.....
جدول ۵-۶- تغییرات کارایی روش SBPSO به تفکیک مراحل و درصد‌های باقیمانده از خصیصه‌ها.....	۶۶.....
جدول ۶-۶- تغییرات کارایی روش MGA به تفکیک مراحل و درصد‌های باقیمانده از خصیصه‌ها.....	۶۹.....
جدول ۷-۶- تغییرات کارایی روش MBPSO به تفکیک مراحل و درصد‌های باقیمانده از خصیصه‌ها.....	۷۲.....
جدول ۱-۷- مقایسه روش SGA و SVM به تفکیک هر بخش از داده‌های LINGSPAM.....	۸۳.....
جدول ۲-۷- مقایسه کارایی سه روش دسته‌بندی هرزنانه.....	۸۴.....
جدول ۳-۷- مقایسه کارایی SGA با فیلترهای غیر SVM.....	۸۵.....
جدول ۴-۷- کارایی روش SGA و روش MBPSO به تفکیک هر یک از ۵ مجموعه.....	۹۰.....
جدول ۱-۸- نتایج آزمون دسته‌بند هر یک از ده کلاس اول رویترز-۲۱۵۷۸ در روش LIU و همکاران.....	۹۸.....
جدول ۲-۸- نتایج آزمون دسته‌بند هر یک از ده کلاس اول رویترز-۲۱۵۷۸ در روش پیشنهادی.....	۹۸.....
جدول ۳-۸- نتایج میانگین‌گیری جزئی حاصل از ده کلاس اول رویترز-۲۱۵۷۸ در روش پیشنهادی.....	۹۸.....
جدول ۴-۸- مقایسه روش پیشنهادی و روش‌های آماری موجود.....	۱۰۰.....

فهرست اشکال

عنوان.....	صفحه.....
شکل ۳-۱- انتخاب زیر مجموعه خصیصه با استفاده از GA و DISTAI.....	۲۴
شکل ۴-۱- الگوریتم اولیه PSO.....	۳۶
شکل ۵-۱- اولین فرایند پیشنهادی انتخاب خصیصه SFSP.....	۴۴
شکل ۵-۲- تابع برازندگی مرحله دوم در GA.....	۴۷
شکل ۵-۳- تابع برازندگی مرحله سوم.....	۴۸
شکل ۵-۴- تابع برازندگی مرحله چهارم.....	۴۹
شکل ۵-۵- تابع برازندگی مرحله سوم.....	۵۲
شکل ۵-۶- تابع برازندگی مرحله چهارم.....	۵۲
شکل ۵-۷- دومین فرایند پیشنهادی انتخاب خصیصه MFSP.....	۵۵
شکل ۵-۸- تابع برازندگی مرحله دوم.....	۵۶
شکل ۶-۱- نمودار تغییرات کارایی روش SGA.....	۶۱
شکل ۶-۲- نمودار تغییرات کارایی روش SGA در فاصله ۰ تا ۲۰.....	۶۱
شکل ۶-۳- نمودار تغییرات دقت در مرحله ۴ روش SGA.....	۶۳
شکل ۶-۴- نمودار تغییرات فراخوانی در مرحله ۴ روش SGA.....	۶۳
شکل ۶-۵- نمودار تغییرات F1 در مرحله ۴ روش SGA.....	۶۳
شکل ۶-۶- نمودار تغییرات کارایی روش SBPSO.....	۶۵
شکل ۶-۷- نمودار تغییرات کارایی روش MGA.....	۶۸
شکل ۶-۸- نمودار تغییرات کارایی روش MGA در فاصله ۰ تا ۲۰.....	۶۸
شکل ۶-۷- نمودار تغییرات کارایی روش MBPSO.....	۷۱
شکل ۶-۸- نمودار تغییرات کارایی روش MBPSO در فاصله ۱/۵ تا ۱۰.....	۷۱
شکل ۶-۱۱- نمودار تغییرات F1 در روش های پیشنهادی انتخاب خصیصه.....	۷۴

- شکل ۶-۱۲- نمودار تغییرات F1 روش‌های انتخاب خصیصه IG، CHI و روش‌های پیشنهادی..... ۷۵
- شکل ۷-۱- تابع برازندگی انتخاب خصیصه جداسازی هرزنامه..... ۸۰
- شکل ۷-۲- مقایسه ارزیابی روش SGA و SVM..... ۸۴
- شکل ۷-۳- ارزیابی روش SGA، TFV+SVM و IG+SVM..... ۸۴
- شکل ۷-۴- نتایج ارزیابی روش SGA، بیزین ساده، SPAMCOP و حییبی و همکارش..... ۸۵
- شکل ۷-۵- نمودار میله‌ای دقت (P) روی هر یک از ۵ مجموعه، قبل و بعد از انتخاب خصیصه..... ۹۱
- شکل ۷-۶- نمودار میله‌ای فراخوانی (R) روی هر یک از ۵ مجموعه، قبل و بعد از انتخاب خصیصه..... ۹۱
- شکل ۷-۷- نمودار میله‌ای صحت روی هر یک از ۵ مجموعه، قبل و بعد از انتخاب خصیصه..... ۹۲
- شکل ۸-۱- مقایسه دقت روش پیشنهادی و دقت روش LIU و همکاران..... ۹۹
- شکل ۸-۲- مقایسه فراخوانی روش پیشنهادی و دقت روش LIU و همکاران..... ۹۹
- شکل ۸-۴- مقایسه روش پیشنهادی و روش‌های آماری موجود..... ۱۰۰

فصل اول

کلیات

۱-۱- مقدمه

امروزه با توجه به گسترش روز افزون اسناد متنی الکترونیکی، توسعه روش‌های پردازش خودکار اسناد اهمیت زیادی پیدا کرده است.

یکی از کارهای بنیادی در مدیریت سندهای متنی، اختصاص اسناد به دسته‌های از پیش تعیین شده می‌باشد که دسته‌بندی (طبقه بندی) متون نامیده می‌شود [۱]. این عمل مستلزم فهم محتوای اسناد و دانش قبلی از حوزه‌های مربوطه است. به همین دلیل، دسته‌بندی اسناد متنی در گذشته توسط افراد متخصص انجام می‌شد. با توجه به این که رویکرد سنتی دسته‌بندی برای حجم زیاد اسناد غیر عملی است، در دهه اخیر تلاش‌های فراوانی در زمینه توسعه الگوریتم‌هایی برای دسته‌بندی خودکار متون صورت گرفته است. در این میان، رویکرد یادگیری ماشین^۱ بسیار مورد توجه بوده است که در آن، با استفاده از مجموعه کوچکی از اسناد برچسب خورده (اسنادی که دسته آنها مشخص است) می‌توان دسته‌بندی را یادگرفت که به صورت خودکار به دسته‌بندی اسناد جدید پردازد [۱]. مهمترین روش‌های استفاده شده در زمینه دسته‌بندی متون عبارتند از: (۱) روش بیزین ساده^۲، (۲) روش Roccio^۳، (۳) روش K همسایه نزدیکتر^۴، (۴) روش‌های رگرسیون، (۵) درخت‌های تصمیم‌گیری، (۶) شبکه‌های عصبی، (۷) روش SVM^۵، (۸) روش‌های بر اساس قاعده^۶ و (۹) روش‌های تکاملی.

به طور کلی، روش‌های دسته‌بندی بر اساس خصیصه‌ها^۶ شکل می‌گیرند و در حوزه متون، خصیصه‌ها همان کلمات هستند. از آنجا که تعداد کلمات در متن بسیار زیاد است، تعداد خصیصه‌ها و در نتیجه تعداد ابعاد فضای خصیصه‌ها بسیار زیاد می‌باشد و این مسئله می‌تواند برای بسیاری از

^۱ Machine Learning

^۲ Naïve Bayesian

^۳ K Nearest Neighbor

^۴ Support Vector Machines

^۵ Decision Rule Classifier

^۶ Features

الگوریتم‌های یادگیری ماشین مشکل‌ساز باشد. در نتیجه، کاهش فضای خصیصه‌ها به صورت خودکار و بدون از دست دادن دقت دسته‌بندی، یکی از اهداف مهم دسته‌بندی متون به شمار می‌آید.

۱-۲- اهداف و نتایج پژوهش

با توجه به بعد بالای فضای خصیصه‌ها در حوزه متون، استفاده از روش‌هایی به منظور حذف خصیصه‌های نامربوط و خصیصه‌هایی که تأثیر چندانی در کارایی دسته‌بندی ندارند، لازم به نظر می‌رسد. همچنین، برخی خصیصه‌ها باعث ایجاد اغتشاش و کاهش کارایی دسته‌بندی می‌شوند. حذف چنین خصیصه‌هایی ضمن کاهش ابعاد فضای خصیصه‌ها، کارایی دسته‌بندی را افزایش می‌دهد.

در این تحقیق با در نظر گرفتن سطح انتظار کاربر از انتخاب خصیصه، دو فرایند انتخاب خصیصه ارائه شده است. فرایند اول، SFSP^۱، شامل ۴ مرحله می‌باشد. در مرحله اول، کارایی دسته‌بند با استفاده از کل خصیصه‌ها و بدون انتخاب خصیصه محاسبه می‌شود. در مرحله دوم، خصیصه‌های نویز یعنی خصیصه‌هایی که حضورشان تأثیری در کارایی دسته‌بندی ندارد و یا آنهایی که باعث کاهش کارایی دسته‌بندی می‌شوند، حذف می‌گردند. به این ترتیب، ضمن کاهش خصیصه‌ها، کارایی دسته‌بند تا حد ممکن بالا می‌رود. در مرحله سوم، سیستم با حفظ کارایی اولیه (نتایج حاصل از مرحله اول) به حذف بیشتر خصیصه‌ها می‌پردازد. در صورتیکه کاربر مایل به کاهش بیشتر تعداد خصیصه‌ها باشد، سیستم وارد مرحله چهارم می‌شود و سعی می‌کند که با دقت کمتر ولی قابل قبولی نسبت به مرحله اول، حذف خصیصه‌ها را ادامه دهد.

به منظور پیاده‌سازی فرایند SFSP ابزاری برای جستجوی فضای خصیصه‌ها و انتخاب زیر مجموعه خصیصه بهینه و نیز ابزاری برای ارزیابی زیر مجموعه خصیصه‌ها نیاز است. از آنجا که پیدا کردن بهترین زیر مجموعه خصیصه یک مسأله NP می‌باشد [۲، ۳]، استفاده از روش‌های حل تقریبی برای یافتن زیر مجموعه خصیصه بهینه مناسب می‌باشد. در نتیجه در این پژوهش، دو روش برای پیاده‌سازی SFSP پیشنهاد شده است که در روش اول، SGA، الگوریتم ژنتیک (GA) [۴] و در روش دوم، SBPSO، نوع دودویی الگوریتم بهینه‌سازی گروهِ ذرات (PSO) [۵، ۶] برای جستجوی

^۱ Sequential Feature Selection Process

^۲ Genetic Algorithm

^۳ Particle Swarm Optimization

زیرمجموعه خصیصه مناسب به کار رفته است. همچنین، از روش دسته‌بندی SVM برای دسته‌بندی متون با توجه به خصیصه‌های منتخب، و ارزیابی زیر مجموعه خصایص استفاده می‌شود.

حضور خصیصه‌های نویز در بردار خصیصه‌ها ممکن است موجب گمراه کردن ابزارهای جستجوی GA یا PSO شود و آنها را از ادامه جستجو برای حذف نویزهای بیشتر باز دارد. برای حل این مشکل، در این تحقیق، فرایند دومی پیشنهاد شده است،¹ MFSP، که در آن، در فواصل زمانی تعیین شده و یا در صورت ارضاء شرط خاصی، نویزهای تعیین شده تا کنون، دور ریخته شده و فضای خصیصه‌ها، به خصیصه‌های انتخاب شده کاهش می‌یابد. سپس حذف نویز در فضای خصیصه‌های جدید، ادامه می‌یابد. این عمل، یعنی کاهش بردار خصیصه‌ها بعد از گذشت چندین تکرار، تا زمانی ادامه می‌یابد که در مرحله جدید قادر به حذف خصیصه‌های بیشتر نباشیم.

برای پیاده سازی MFSP دو روش پیشنهاد شده است که در روش اول، MGA، از الگوریتم ژنتیک و در روش دوم، MBPSO، از نوع دودویی الگوریتم بهینه سازی گروه ذرات برای جستجوی زیرمجموعه خصیصه مناسب به کار رفته است. در این دو روش نیز، از روش دسته‌بندی SVM برای دسته‌بندی متون با توجه به خصیصه‌های منتخب، و ارزیابی زیر مجموعه خصایص استفاده می‌شود.

نتایج ارزیابی‌ها نشان می‌دهد، در کل می‌توان نتیجه گرفت، در صورتی که کم بودن تعداد خصیصه‌های انتخابی و بالا بودن معیار F1 از اهمیت یکسانی برخوردار باشند و همچنین در شرایطی که بالا بودن معیار F1 اهمیت بیشتری داشته باشد، روش MBPSO بهترین روش بوده و روش‌های MGA، SGA و SBPSO به ترتیب در رده‌های بعدی قرار دارند. در صورتی که کمتر بودن تعداد خصیصه‌های انتخاب شده از اهمیت بیشتری نسبت به بالا بودن معیار F1 برخوردار باشد، روش MGA بهترین بوده و سپس روش‌های MBPSO، SGA و SBPSO به ترتیب در رده‌های بعدی قرار دارند.

همچنین مقایسه روش‌های SGA و SBPSO و همچنین روش‌های MGA و MBPSO، نشان می‌دهد که GA و BPSO، هر دو در قالب MFSP بهتر از SFSP عمل می‌کنند و میزان بهبود عملکرد BPSO در MFSP، بیشتر از بهبود عملکرد GA در SFSP مشهود است.

¹ Multi start Feature Selection Process

مقایسه روش‌های پیشنهادی با روش‌های انتخاب خصیصه IG و CHI که دو تا از بهترین روش‌های انتخاب خصیصه فیلتری در حوزه متن می‌باشند، نشان داد که روش‌های پیشنهادی، نتایج بهتری را نسبت به روش‌های انتخاب خصیصه IG و CHI از خود نشان می‌دهند. در ضمن، روش IG بعد از حذف حدود ۸۵ درصد خصیصه‌ها و روش CHI بعد از حذف حدود ۹۰ درصد خصیصه‌ها، دیگر قادر به حفظ کارایی اولیه نمی‌باشند.

به منظور سنجش کارایی روش‌های پیشنهادی در سایر دامنه‌ها، روش SGA را در راستای انتخاب خصیصه در جداسازی هرزنامه‌های متنی بکار برده‌ایم. نتایج اعمال این روش روی مجموعه استاندارد [7] LingSpam نشان داده است که روش پیشنهادی علاوه بر جداسازی هرزنامه‌ها با دقت ۱ و فراخوانی ۰/۹۹۸، قدرت کاهش خصایص تا حدود ۱۲/۷۴٪ تعداد اولیه را دارد. همچنین، مقایسه روش پیشنهادی با سایر روش‌های جداسازی هرزنامه‌ها، نشان داده است که روش پیشنهادی به خوبی قابل مقایسه با بهترین روش‌های آماری موجود است. همچنین روش‌های SGA و MBPSO را به منظور انتخاب حداقل ژن‌های لازم برای دسته‌بندی درست نسوج تومور به کار برده‌ایم. روش‌های پیشنهادی بر روی مجموعه داده سرطان روده بزرگ آزمایش شده‌اند. نتایج نشان داده‌اند که روش‌های SGA و MBPSO برای انتخاب ژن و دسته‌بندی نسوج تومور مناسب بوده و می‌توانند با انتخاب تعداد بسیار کم (به طور متوسط ۳) ژن، کارایی دسته‌بندی را افزایش دهند. همچنین روش‌های پیشنهادی می‌توانند با شناساندن ژن‌های مهم در دسته‌بندی نسوج تومور در سایر جنبه‌های پزشکی و درمان مفید واقع شوند.

آخرین کار انجام شده در این پژوهش، استفاده از روش‌های تکاملی برای یادگیری دسته‌بند متون است. به این منظور، روش Liu و همکارانش [۸] به عنوان یکی از روش‌های مستقل از ساختار متن بر روی داده‌های استاندارد روترز آزمون شده و با روش‌های آماری مقایسه شده است. در این پژوهش، تابع ارزیابی جدیدی تعریف شده است که موجب بهبود کارایی روش Liu و همکارانش، از نظر دقت: ۲٪، از نظر فراخوانی ۱۵٪، از نظر FI ۹٪ و از نظر صحت ۲٪ بهتر از روش Liu و همکارانش عمل می‌کند. البته روش بهبود یافته دسته‌بندی نسبت به روش‌های آماری موجود ضعیفتر می‌باشد.

۱-۳- مروری بر ساختار پایان نامه

این پایان نامه شامل ۹ فصل است. فصل دوم به بررسی تاریخیچه پژوهش در زمینه انتخاب خصیصه و دسته‌بندی متون پرداخته و روش‌های مطرح آماری و همچنین روش‌های برگرفته از طبیعت در زمینه انتخاب خصیصه در حوزه متن و غیر متن را شرح می‌دهد. در نهایت نیز انواع روش‌های مطرح دسته‌بندی متون توضیح داده می‌شود.

فصل سوم درباره متون و مفاهیم پایه مربوط به آن از جمله خصیصه، روش‌های نمایش اسناد متنی، مجموعه‌های آموزش و آزمون، مجموعه داده‌های استاندارد و معیارهای ارزیابی کارایی دسته‌بندی متون، صحبت می‌کند. در نهایت نیز، روش دسته‌بندی SVM به دلیل استفاده از آن در این پژوهش، به صورت دقیقتر توضیح داده می‌شود.

فصل چهارم به توصیف روش‌های محاسباتی برگرفته از طبیعت از جمله الگوریتم‌های تکاملی و روش بهینه‌سازی گروه ذرات، انواع آنها و مفاهیم مربوطه می‌پردازد.

در فصل پنجم دو فرایند پیشنهادی انتخاب خصیصه توضیح داده شده و هر یک از آنها با استفاده از دو روش GA و PSO پیاده‌سازی می‌شوند.

نتایج حاصل از آزمون روش‌های پیشنهادی انتخاب خصیصه روی داده‌های روترز، در فصل ششم آورده شده است. همچنین روش‌های پیشنهادی با بهترین روش‌های آماری انتخاب خصیصه متون مقایسه می‌شوند.

فصل هفتم به بررسی شایستگی روش‌های پیشنهادی در سایر دامنه‌ها اختصاص دارد. بخش اول فصل هفتم به بحث جداسازی هرزنامه‌های متنی اختصاص دارد. در این بخش، ابتدا مقدمه‌ای بر هرزنامه‌ها و روش‌های جداسازی آنها گفته می‌شود. سپس پژوهش‌های مرتبط در زمینه روش‌های جداسازی هرزنامه مطرح می‌شود. در ادامه، تابع ارزیابی و سایر شرایط محیطی استفاده شده برای مسأله انتخاب خصیصه در حوزه هرزنامه‌های متنی، توضیح داده می‌شود. سپس، نتایج ارزیابی کارایی روش SGA در مقایسه با روش‌های موجود مطرح می‌گردد.

بخش دوم فصل هفتم نیز به بحث دسته‌بندی نسوج تومور اختصاص دارد. ابتدا مقدمه‌ای بر دسته‌بندی نسوج تومور و روش‌های نوین در این زمینه گفته می‌شود. سپس پژوهش‌های مرتبط در زمینه دسته‌بندی نسوج تومور مطرح می‌گردد. در ادامه، تابع ارزیابی و سایر شرایط محیطی استفاده